

第一章 预 篇

§1 引 言

数值代数是计算数学的重要分支，它是研究各类代数问题的数值解法的理论，例如矩阵的变换和分解、线性和非线性代数方程的解法、矩阵的特征值及特征向量的求解问题等。

量子力学、结构力学、网络分析、大地测量、机械振动、管理科学等应用技术学科中提出的大量的计算问题，大都直接归属于计算数学。而计算数学各分支中的问题，如有限元方法、谱分解理论、稀疏矩阵问题、最小二乘法问题、数学规划与组合数学问题又都直接或间接成为数值代数的重要源泉和广阔的应用场所。因此，数值代数的研究和发展不仅直接推动着计算数学的发展而且也是其它技术学科发展的重要前提。

古典代数理论，对于代数问题解的研究，已日臻完善；但当未知数的个数较多时，其数值解的获得，则仍是一个困难的问题。即使在快速电子计算机问世以来，这个问题所面临的挑战，也并未因此而减弱。相反，随着科学技术的发展，当生产斗争、科学实验日益迫切地要求人类向高维问题进军的时候，新的、富有成效的、收敛性能和稳定性能都更为良好的各种算法的研制，更是迫切地提到日程上来。和其它计算数学分支一样，数值代数也随着近代物质文明的需要而日新月异地改变着自己的面貌。

本书作为数值分析的后继课程，将结合数值代数的近代发

展，把数值分析中有关数值代数的章节，适当地给于扩充和引申，同时也为一些新的实用的方法的学习和探索提供初步的导引。对于重要的方法，将以接近语句的形式把它们写出来。自然，算法描述是理论分析和推理演算的结晶，也是计算的有效性和经济性的统一，而不是规则和条文的堆砌。

本章将讨论一些在以后各章节中都要用到的工具和定理，它们包括：初等矩阵、矩阵的分解以及有关特征扰动方面的问题。

今将全书统一使用的符号一并在此列出。除维数和足码以外，通常用小写希腊字母代表数，用小写黑体拉丁字母代表向量，用大写的希腊或拉丁字母代表矩阵； A^T 表示矩阵 A 的转置， A^H 则表示它的共轭转置；用 R^n 、 C^n 分别表示 n 维实的和复的线性空间、用 $R^{n \times m}$ 、 $C^{n \times m}$ 分别表示 n 行 m 列的实或复的矩阵所成的空间；矩阵 A 的列一般记为 a_i ，而以 a_{ij} 或 a_{ji} 表示其在位置 (i, j) 的元素；用 e_i 表示单位矩阵 I 的第 i 列，而

$$e = \sum e_i$$

则表示分量全为 1 的向量； $|x|$ 、 $|A|$ 分别表示用向量 x 和矩阵 A 的分量的绝对值做成的向量和矩阵。在不致发生误解时，单维阵的维数、零阵或零向量的维数将不特别指明。用 $\det(A)$ 、 $\lambda(A)$ 分别表示矩阵 A 的行列式和它的一个特征值，而用

$$\text{tr}(A) = \sum_{i=1}^n a_{ii} \quad (1.1.1)$$

$$\text{和} \quad \rho(A) = \max |\lambda(A)| \quad (1.1.2)$$

分别表示 A 的迹和谱半径；用

$$A \begin{pmatrix} i_1, i_2, \dots, i_p \\ j_1, j_2, \dots, j_p \end{pmatrix} \quad (1.1.3)$$

表示选自 A 的第 i_1, \dots, i_p 各行和 j_1, \dots, j_p 各列的交点处的元素所成的行列式。

在不作特殊解释时, 向量范数 $\|\mathbf{x}\|$ 恒指 l_2 空间中的范数, 即

$$\|\mathbf{x}\| = (\mathbf{x}^H \mathbf{x})^{1/2}, \quad (1.1.4)$$

而矩阵范数 $\|A\|$ 是指与向量范数相容或从属于它的一种范数, 即

$$\|A\mathbf{x}\| \leq \|A\| \|\mathbf{x}\| \quad (1.1.5)$$

或再加上一个限制: (1.1.5) 中的等号至少对一个 $\mathbf{x} \neq 0$ 成立.

§2 初等矩阵

我们把形如

$$E(\mathbf{u}, \mathbf{v}, \sigma) = I - \sigma \mathbf{u} \mathbf{v}^H \quad (1.2.1)$$

的矩阵叫作**初等阵**, 其中 \mathbf{u} 、 \mathbf{v} 是 n -维向量, σ 是数, 该阵是一个单位阵与一个秩 1 矩阵的组合. 当要把一个一般形式的矩阵约化为某种标准形式或任何压缩形式时, 常需对矩阵进行一系列简单的相似变换, 而这种工作通常是借助初等阵来实现的.

首先指出在 (1.2.1) 中, 当视 σ 为参数, \mathbf{u} 、 \mathbf{v} 为固定向量时, 矩阵类 (1.2.1) 对乘法是封闭的. 实际上

$$E(\mathbf{u}, \mathbf{v}, \sigma) E(\mathbf{u}, \mathbf{v}, \tau) = E(\mathbf{u}, \mathbf{v}, \lambda), \quad (1.2.2)$$

其中,

$$\lambda = \sigma \tau \mathbf{v}^H \mathbf{u} - \sigma - \tau. \quad (1.2.3)$$

由于 $E(\mathbf{u}, \mathbf{v}, 0) = I$, 令 $\lambda = 0$ 时解得

$$\tau = -\sigma(1 - \sigma \mathbf{v}^H \mathbf{u}), \quad 1 - \sigma \mathbf{v}^H \mathbf{u} \neq 0, \quad (1.2.4)$$

(1.2.4) 同时给出了 $E(\mathbf{u}, \mathbf{v}, \sigma)$ 有逆的条件和逆阵中的 τ 的计算.

直接计算 $E(u, v, \sigma)$ 的行列式, 还可进一步证明

$$\det E(u, v, \sigma) = 1 - \sigma v^H u. \quad (1.2.5)$$

在消元中常见的有如下的初等阵.

2.1 行交换矩阵

$$P_{ij} = E(e_i - e_j, e_i - e_j, 1)$$

除第 i 行、 j 行、 i 列及 j 列上的 4 个元素有如下分布外, 其余元素与单位矩阵 I 相同, 称为行交换矩阵.

$$\begin{array}{cc} i \text{ 列} & j \text{ 列} \\ 0 & 1 & i \text{ 行} \\ 1 & 0 & j \text{ 行} \end{array}$$

由于: $P_{ij}P_{ij} = I$ 及 P_{ij} 实对称, 故 P_{ij} 是简单的正交阵. 当用 P_{ij} 左乘某一矩阵时, 其结果将使原矩阵的 i 、 j 两行互换, 右乘时则将使 i 、 j 两列互换.

2.2 排列矩阵

$$P_{a_1, \dots, a_n}$$

称为排列矩阵, 其中 (a_1, \dots, a_n) 代表 $(1, 2, \dots, n)$ 的一个排列, 其元素, 除位于

$$(i, a_i), i = 1, 2, \dots, n$$

的元素规定为 1 外, 其余全为零. 该矩阵在每个行和列上都恰有一个 1, 因此, 其转置阵也是排列阵, 并且

$$P_{a_1, \dots, a_n}^{-1} = P_{a_1, \dots, a_n}^T, \quad (1.2.6)$$

因而是正交阵, 排列阵均可表示成 P_{ij} 型矩阵的乘积. 它实际是单位阵经有限次换行的结果, 因此其行列式的值均为 ± 1 . 据定义, 有

$$P_{a_1, \dots, a_n} = \sum_{i=1}^n e_{a_i} e_i^T. \quad (1.2.7)$$

当用 P_{a_1, \dots, a_n} 左乘某矩阵时, 其结果将把原矩阵的第 a_i 行换到

第 i 行, 右乘时, 则结果是把原矩阵的第 α_i 列换到第 i 列。

2.3 消元矩阵

形如

$$\begin{aligned} L_i &= E(l_i, e_i, 1), \\ e_k^T l_i &= 0, \quad k \leq i \end{aligned} \quad (1.2.8)$$

的矩阵, 其中 l_i 是前 i 个分量为零的向量, 设

$$l_i^T = (0, \dots, 0, l_{i+1,i}, \dots, l_{n,i}),$$

则

$$L_i = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & -l_{i+1,i} & 1 & \\ & & \vdots & & \ddots \\ & & -l_{n,i} & & & 1 \end{pmatrix} \quad (1.2.9)$$

不难验证

$$L_i^{-1} = E(-l_i, e_i, 1), \quad \det(L_i) = 1. \quad (1.2.10)$$

当用 L_i 左乘于 A 时, 其结果将把原矩阵中的第 i 行分别乘以 $-l_{i+1,i}, \dots, -l_{n,i}$ 并依次加到矩阵的第 $i+1, \dots, n$ 各行中, 右乘时, A 的第 $i+1$ 至第 n 列将分别乘以 $-l_{i+1,i}, -l_{i+2,i}, \dots, -l_{n,i}$ 并加到第 i 列, 即左乘时将改变原矩阵的第 i 行以下的所有的行, 右乘时只改变原矩阵的第 i 列。这种矩阵在求逆时, 仅需把 l_i 换为 $-l_i$, 实际并不进行算术运算, 因而并不引进舍入误差。不仅如此, 在计算形如

$$L_1 L_2 \cdots L_k, \quad L_1^{-1} L_2^{-1} \cdots L_k^{-1}$$

的矩阵乘积时, 也不会引进误差。例如乘积矩阵 $L_1 \cdots L_k$ 的第 j 列只不过是 L_i 的第 j 列的重抄而已。类似地, 若记

$$R_i = L_i^T = E(e_i, l_i, 1),$$

可知

$$R_k^{-1} R_{k-1}^{-1} \cdots R_1^{-1} = \begin{pmatrix} 1 & x & \cdots & x \\ & 1 & x & \cdots & x \\ & & \ddots & & \\ & & & 1 & \cdots & x \\ 0 & & & & \ddots & 1 \end{pmatrix}$$

应是单位上三角阵, 其第 1 至 k 行将依次是 $l_1^T \cdots l_k^T$ 的重抄.

2.4 初等 Hermite 矩阵

$$\text{矩阵 } H = E(u, u, 2), \quad \|u\| = 1 \quad (1.2.11)$$

称为初等 Hermite 矩阵, 文献上由于强调的侧面不同, 有时又被称为 Householder 矩阵或镜像反射矩阵.

不难验证它既是 Hermite 阵又是酉阵. 由于它还是一个对合阵, 因而它的逆就是其自身. 据 (1.2.5) 及 (1.2.11) 易知, $\det(H) = -1$.

由于下面三个定理, 使这一类矩阵在矩阵的变换与分解中有着十分重要的应用.

定理 1.2.1 对任意 $x \in C^n$, 存在初等 Hermite 阵 H , 使

$$Hx = \alpha e_1, \quad (1.2.12)$$

其中 $|\alpha| = \|x\|$; 特别, 可取

$$\alpha = -\text{sign}(\text{Re}(\xi_1)) \|x\|, \quad (1.2.13)$$

其中 ξ_1 是 $x^T = (\xi_1, \cdots, \xi_n)$ 的第 1 个分量.

证明 记 $H = I - 2uu^H$, 由

$$Hx = x - 2(u^H x)u = \alpha e_1,$$

故

$$u = \frac{(x - \alpha e_1)}{\|x - \alpha e_1\|} e^{i\theta} \quad (\theta \text{ 是实数}),$$

但 uu^H 与 θ 无关, 取 $\theta = 0$. 考查

$$\|x - \alpha e_1\|^2 = 2(\|x\|^2 - \text{Re}(\alpha \xi_1)). \quad (1.2.14)$$

当按 (1.2.13) 确定 α 时, (1.2.14) 右端中才能有更多的有效数字.

作为本定理的应用, 可证明下述著名的 **Schur** 定理.

定理 1.2.2 任意矩阵可以经由酉相似变换化为上三角阵.

推论 1.2.1 正规矩阵可以经由酉相似变换化为对角阵.

由于定理 1.2.2 及其推论的证明, 并借助矩阵的特征值及特征向量, 可知, 对计算矩阵的特征问题来说, 定理本身并非是构造性的, 但它却是构造逐次逼近法求解矩阵特征问题的基础.

定理 1.2.2 及其推论的证明, 已收入本章的习题. 在第三章中, 我们就实矩阵的情形, 给出了一个类似的更为适用的结果.

§3 矩阵的分解

首先, 我们指出: 对于任意的矩阵 $A \in \mathbb{C}^{n \times m}$, 设其秩为 p , 则可以找到秩为 p 的两个矩阵 M 和 P , 其中 M 有 p 个列, P 有 p 个行, 即 $M \in \mathbb{C}^{n \times p}$ 而 $P \in \mathbb{C}^{p \times m}$, 使

$$A = MP. \quad (1.3.1)$$

我们将用两种方法说明上述论断.

3.1 列主元法

$$\text{设 } A = A^{(0)} = (a_1^{(0)}, a_2^{(0)}, \dots, a_m^{(0)}) \neq 0. \quad (1.3.2)$$

令 j_1 为使 $a_{j_1}^{(0)} \neq 0$ 的最小的足码, 且令

$$|e_{j_1}^T a_{j_1}^{(0)}| = \max_i |e_i^T a_{j_1}^{(0)}|, \quad (1.3.3)$$

即 $e_{j_1}^T a_{j_1}^{(0)}$ 是 $a_{j_1}^{(0)}$ 的依模最大的分量, 作

$$\tilde{A}^{(1)} + P_{1j_1} A^{(0)} = (\tilde{a}_1^{(0)}, \tilde{a}_2^{(0)}, \dots, \tilde{a}_m^{(0)}), \quad (1.3.4)$$

记

$$l_1^T = (1, l_{21}, \dots, l_{n1}),$$

$$l_{k1} = e_k^T \tilde{a}_{j_1}^{(0)} / e_{j_1}^T \tilde{a}_{j_1}^{(0)}, \quad k = 2, \dots, n, \quad (1.3.5)$$

$$L_1 = E(l_1, e_1, 1),$$

$$A^{(1)} = L_1 \tilde{A}^{(1)},$$

则在

$$A^{(1)} = (a_1^{(1)}, a_2^{(1)}, \dots, a_m^{(1)}) \quad (1.3.6)$$

中, 其第一行因 $e_1^T a_{j_1}^{(1)} \neq 0$ 必为非零行, 但其前 $j_1 - 1$ 列上的元素以及 $a_{j_1}^{(1)}$ 的第一行以下的元素全为0.

如果 $A^{(1)}$ 的第一行以下全为零, 则由

$$A^{(0)} P_{1j_1} L_1^{-1} A^{(1)} = P_{1j_1} L_1^{-1} A^{(1)} \quad (1.3.7)$$

及 $A_1 = e_1 e_1^T A_1$, 有

$$L_1^{-1} A^{(1)} = (L_1^{-1} e_1) (e_1^T A^{(1)}).$$

可见 $A^{(0)}$ 已经表示成1个非零的列向量 $P_{1j_1} L_1^{-1} e_1$ 与一个非零的行向量 $e_1^T A^{(1)}$ 的乘积, 此时 $A^{(0)}$ 的秩也只能是1. 即所述论断是正确的.

如果 $A^{(1)}$ 的第1行以下第 j_1 列以右的矩阵 $\tilde{A}^{(2)}$ 非零, 在 $A^{(0)}$ 的秩为 p 的前提下可知 $\tilde{A}^{(2)}$ 的秩应为 $p - 1$. 对行 n 使用归纳法, 可设

$$\begin{aligned} A^{(1)} &= (A_1^{(1)} : A_2^{(1)}) = \begin{pmatrix} 0 \cdots a_{j_1}^{(1)} & a_{j_1+1}^{(1)} \cdots a_m^{(1)} \\ \cdots & \cdots \\ 0 & \tilde{A}_2 \end{pmatrix} \\ &= \begin{pmatrix} A_{1,1}^{(1)} : A_{1,2}^{(1)} \\ \cdots \\ 0 : \tilde{A}_2 \end{pmatrix}, \end{aligned}$$

而据归纳法可设 $\tilde{A}_2 = \tilde{M} \tilde{P}$, 这里 \tilde{M} , \tilde{P} 都是秩为 $p - 1$ 且各有 $p - 1$ 列和 $p - 1$ 行的阵, 记 $\bar{P} = (0 : \tilde{P})$, \bar{P} 是秩为 $p - 1$ 但有 $n - 1$ 行和 m 列的阵. 可得

$$A^{(1)} = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} (A_{1,1}^{(1)} \cdots A_{1,2}^{(1)}) + \begin{pmatrix} 0 \\ \vdots \\ \tilde{M} \end{pmatrix} \tilde{P}$$

$$= \begin{pmatrix} 1 & \vdots & 0 & A_{1,1}^{(1)} & \vdots & A_{1,2}^{(1)} \\ 0 & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \vdots & \tilde{M} & 0 & \vdots & \tilde{P} \end{pmatrix},$$

右端已约化为所需要的形式。由于

$$A^{(0)} = P_{1:n} L_1^{-1} \begin{pmatrix} 1 & & A_{1,1}^{(1)} & A_{1,2}^{(1)} \\ & 0 & & \\ & \vdots & \tilde{M} & \\ 0 & & 0 & \tilde{P} \end{pmatrix},$$

令

$$P_{1:n} L_1^{-1} \begin{pmatrix} 1 & & & \\ & 0 & & \\ & \vdots & \tilde{M} & \\ & 0 & & \end{pmatrix} = M, \quad \begin{bmatrix} A_{1,1}^{(1)} & A_{1,2}^{(1)} \\ 0 & \tilde{P} \end{bmatrix} = P,$$

可知

$$A^{(0)} = MP \quad (1.3.8)$$

具有所要求的形式。

在 $n=m$ 、 $A^{(0)}$ 非奇异的前提下，如果进一步假定 $M=L$ 是单位下三角阵， $P=R$ 是上三角阵，此时，称 $A=A^{(0)}$ 有三角分解：

$$A = LR. \quad (1.3.9)$$

可证此种分解是唯一的。

3.2 初等 Hermite 阵法

这里，我们主要有以下结论。

设 $A \in \mathbb{C}^{n \times n}$ 且秩为 p ，则存在秩为 p 的两个矩阵 W 和 R ，它们各有 p 个列和 p 个行，其中 W 的列还是互相直交的单位向量，使

$$A = WR, \quad (1.3.10)$$

记

$$A = A^{(0)} = (a^{(0)}_1, \dots, a^{(0)}_m), \quad (1.3.11)$$

令 $\alpha_{1,0}^{(0)}$ 是 $A^{(0)}$ 中的第一个非零列, 依定理 1.2.1 作初等 Hermite 阵

$$H_1 = I - 2u_1 u_1^H \quad (1.3.12)$$

使

$$H_1 \mathbf{a}_i^{(0)} = \gamma_{1i} \mathbf{e}_1, \quad |\gamma_{1i}| = \|\mathbf{a}_i^{(0)}\| \neq 0, \quad (1.3.13)$$

◆

$$A^{(1)} = H_1 A^{(0)}, \quad (1.3.14)$$

此时, $A^{(i)}$ 的前 $i_1 - 1$ 列是零向量, 第 i_1 列除第 1 个分量 $\gamma_{1i_1} \neq 0$ 外, 其余分量都是零.

如果 $A^{(1)}$ 的第 1 行以下全为零, 这时, 由

$$\mathbf{A}^{(0)} = \mathbf{H}_1^{-1} \mathbf{A}^{(1)} = \mathbf{H}_1 \mathbf{A}^{(1)} = (\mathbf{H}_1 \mathbf{e}_1) (\mathbf{e}_1^T \mathbf{A}^{(1)}) \quad (1.3.15)$$

以及 $H_1 e_1$ 是非零列向量, $e_1^T A^{(1)}$ 是非零行向量和 $A^{(0)}$ 的秩是 1, 故已合于预期的结论.

如果 $A^{(1)}$ 的第 1 行以下第 j_1 列以右的矩阵 $\bar{A}^{(1)}$ 非零, 用归纳法, 可设

$$\bar{A}^{(2)} = \hat{H}^{(2)} \bar{R}^{(2)}. \quad (1.3.16)$$

在 $A^{(0)}$ 的秩为 p 的前提下, $\bar{A}^{(1)}$ 、 $\hat{H}^{(2)}$ 、 $\tilde{R}^{(2)}$ 的秩均应是 $p-1$, 且 $\hat{H}^{(2)}$ 由 $p-1$ 个互相直交的单位向量组成, $\tilde{R}^{(2)}$ 由 $p-1$ 个线性无关的行向量组成. 这时, 由

$$A^{(1)} = \begin{pmatrix} 0 \cdots a_{1,j_1}^{(1)} & : & a_{j_1+1}^{(1)} \cdots a_{l_m}^{(1)} \\ \vdots & & \vdots \\ 0 & : & \tilde{H}^{(2)} \tilde{R}^{(2)} \end{pmatrix} = \begin{pmatrix} A_{l,1}^{(1)} & : & A_{l,2}^{(1)} \\ \vdots & & \vdots \\ 0 & : & \tilde{H}^{(2)} \tilde{R}^{(2)} \end{pmatrix}$$
$$= \begin{pmatrix} A_{l,1}^{(1)} & A_{l,2}^{(1)} \\ \vdots & \vdots \\ 0 & \end{pmatrix} + \begin{pmatrix} 0 \\ \vdots \\ \tilde{H}^{(2)} \tilde{R}^{(2)} \end{pmatrix}.$$

这里

$$\bar{R}^{(2)} = (\overset{j_1}{0} \vdots \hat{R}^{(2)})_{(p-1)},$$

故

$$A^{(1)} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \tilde{H}^{(2)} & \ddots & \vdots \\ 0 & & & \end{pmatrix} \begin{pmatrix} A_{1,1}^{(1)} & A_{1,2}^{(1)} \\ & \ddots & \ddots \\ 0 & & \tilde{R}^{(2)} \end{pmatrix},$$

令

$$W = H_1 \begin{pmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \tilde{H}^{(2)} & \ddots & \vdots \\ 0 & & & \end{pmatrix}, \quad R = \begin{pmatrix} A_{1,1}^{(1)} & A_{1,2}^{(1)} \\ & \ddots & \ddots \\ 0 & & \tilde{R}^{(2)} \end{pmatrix},$$

则由 $A^{(0)} = H_1 A^{(1)}$, 有

$$A^{(0)} = A = WR \quad (1.3.17)$$

具有所要求的形式.

当 $n=m$, 且 $A^{(0)}$ 非奇异时, 如果进一步限制 R 为对角元为正数的上三角阵 U , 则分解是存在且唯一的.

首先, 当 A 非奇异时, 由于 $A=WR$ 中, W 、 R 的秩均是 n , 又 W 的各列是互相正交的单位向量, 此时 W 是酉阵, R 可以写成三角阵的形式, 它的对角元一定可以表示成正数, 因为 R 的对角元为

$$\gamma_{kk} = |\gamma_{kk}| e^{i\theta_k}, \quad \theta_k \text{ 是实数} \quad (1.3.18)$$

$$k = 1, 2, \dots, n$$

由

$$WR = (WD)(D^{-1}R) = QU \quad (1.3.19)$$

这里, $D = \text{diag}(e^{i\theta_1}, \dots, e^{i\theta_n})$, 记 $D^{-1}R = U$ 是对角元为正数的上三角阵, 又 D 是对角酉阵, 于是 $Q = WD^{-1}$ 仍是酉阵. 从而

$$A = QU, \quad (1.3.20)$$

其中 Q 是酉阵, U 是对角元为正数的上三角阵. 分解的唯一性, 则可由下面的分析说明, 设

$$Q_1 U_1 = Q_2 U_2,$$

则

$$Q_2^H Q_1 = U_2 U_1^{-1}$$

等式两端既是酉阵又是单位上三角阵，推知两端应同时是单位阵，于是 $Q_1 = Q_2$, $U_1 = U_2$.

3.3 三角分解的表现

当 $A \in \mathbb{C}^{n \times n}$ 非奇异时，有可能得到下列的分解

$$A = LR,$$

其中 L 是单位下三角阵， R 是对角元非零的上三角阵，如果把 R 的对角元从 R 的左边提出来，可得

$$A = LDU, \quad (1.3.21)$$

其中 L 是单位下三角阵， U 是单位上三角阵， D 是非奇异对角阵。下面讨论怎样用 A 的元素来表示 L 、 D 、 U 。

利用矩阵的分块表示法，设

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} D_{11} & \\ & D_{22} \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix}, \quad (1.3.22)$$

这里，再设 A_{11} , L_{11} , U_{11} 都是 $p (\leq n)$ 阶方阵，有

$$A_{11} = L_{11} D_{11} U_{11},$$

$$\det(A_{11}) = \det(D_{11}) = \delta_1 \delta_2 \cdots \delta_p, \quad (1.3.23)$$

$$D_{11} = \text{diag}(\delta_1, \cdots, \delta_p),$$

$$\delta_p = A \begin{pmatrix} 1, 2, \cdots, p \\ 1, 2, \cdots, p \end{pmatrix} / A \begin{pmatrix} 1, 2, \cdots, p-1 \\ 1, 2, \cdots, p-1 \end{pmatrix}, \quad (1.3.24)$$

由

$$(A_{11}, A_{12}) = L_{11}(D_{11}U_{11}, D_{11}U_{12}); \quad (1.3.25)$$

由于两端都是行数为 p 的矩阵，故当两端取列号相同的 p 个列时，得出的行列式也应相等。特别，当两端取 $1 \cdots (p-1)$ 列并附加以第 j 列 ($j > p$) 时，得

$$A \begin{bmatrix} 1, 2, \dots, p-1, p \\ 1, 2, \dots, p-1, j \end{bmatrix} = \det(L_{11}) \det(D_{11}) U \begin{bmatrix} 1, 2, \dots, p-1, p \\ 1, 2, \dots, p-1, j \end{bmatrix},$$

$$u_{p,j} = A \begin{bmatrix} 1, 2, \dots, p-1, p \\ 1, 2, \dots, p-1, j \end{bmatrix} / A \begin{bmatrix} 1, 2, \dots, p \\ 1, 2, \dots, p \end{bmatrix}, \quad (1.3.26)$$

这里, $u_{p,j}$ 是 U 的位于 (p, j) 的元素. 比较

$$\begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} = \begin{bmatrix} L_{11} & D_{11} & U_{11} \\ L_{21} & D_{22} & U_{11} \end{bmatrix}$$

两端各有 n 个行和 p 个列, 两端取前 $p-1$ 行和第 i 行形成行列式, 由于 $\det(U_{11}) = 1$, 得

$$l_{i,j} = A \begin{pmatrix} 1, 2, \dots, p-1, i \\ 1, 2, \dots, p-1, p \end{pmatrix} / A \begin{pmatrix} 1, \dots, p \\ 1, \dots, p \end{pmatrix} \quad (i > p). \quad (1.3.27)$$

§4 矩阵的 Jordan 标准型

设 $A \in \mathbb{C}^{n \times n}$, 据线性代数知识, 存在非奇异矩阵 H , 使

$$H^{-1}AH = J = \begin{pmatrix} J_1 & & \\ & J_2 & \\ & & \ddots \\ & & & J_r \end{pmatrix} \quad (1.4.1)$$

其中

$$J_i = \begin{pmatrix} \lambda_i & 1 & & \\ & \lambda_i & 1 & \\ & & \ddots & \\ & & & \lambda_i & 1 \\ & & & & \lambda_i \end{pmatrix}, \quad (1.4.2)$$

$\lambda_i (i = 1, \dots, r)$ 是 A 的特征值, J_i 为对应 λ_i 的上 Jordan 块, λ_i 可能相异或相等. 为确定起见, 我们不妨规定: J 中对应同一特征值的 Jordan 块是彼此相邻的. 于是, 当用 $J^{(k)}$ 的上标 k 表示 A 的相异特征值的编号, 下标 s 表示对应同一特征值的 Jordan 块的编号时, 则 (1.4.1) 可重写为

$$H^{-1}AH = \text{diag}(J_1^{(1)}, \dots, J_1^{(1)}, \dots, J_r^{(r)}, \dots, J_r^{(r)}). \quad (1.4.3)$$

我们把 (1.4.1) 和 (1.4.3) 统称为 A 的 **Jordan 标准型**, 而称 H 为 A 的一个 **变换矩阵**.

记 $J_s^{(k)}$ 的阶数为 $n_s^{(k)}$, 则

$$n = \sum_{k=1}^r \sum_{s=1}^{i_k} n_s^{(k)}. \quad (1.4.4)$$

当每个 $n_s^{(k)}$ 均为 1 时, A 将相似于一个对角阵, 即

$$H^{-1}AH = \text{diag}(\lambda_1, \dots, \lambda_n). \quad (1.4.5)$$

据习题 (1.12), 任意正规阵均酉相似于一对角阵, 于 (1.4.5) 两端各用 U^H 和 U 去左乘和右乘, 这里 U 是酉阵, 则于右端得到正规阵, 因此, 凡相似于对角阵的阵, 必也相似于一正规阵, 所以, 相似于对角阵的阵又被称为可正规化矩阵, 但酉相似于对角阵的阵则必是正规阵.

一般情况下, 由 (1.4.1) 有

$$AH = HJ = H \text{diag}(J_1, \dots, J_r);$$

记 p 为 J_s ($1 \leq s \leq r$) 的第 1 列在 J 中的编号, 则由

$$AH\mathbf{e}_p = A\mathbf{h}_p = \lambda_s \mathbf{h}_p$$

这里 \mathbf{h}_p 是 H 的第 p 列. 上述表明, 它是 A 的对应于 λ_s 的特征向量. 当 J_s 的阶数大于 1 时, H 的对应于 J_s 的其余的列, 例如对应 J_s 的第 2, 第 3... 列的列, 由于满足

$$A\mathbf{h}_{p+1} = \lambda_s \mathbf{h}_{p+1} + \mathbf{h}_p, A\mathbf{h}_{p+2} = \lambda_s \mathbf{h}_{p+2} + \mathbf{h}_{p+1}, \dots \quad (1.4.7)$$

因而都不是 A 的特征向量. 据此, 可知 A 的线性无关的特征向量的个数应为 A 的 Jordan 标准型中的 Jordan 块的个数, 即 r .

于是, 一般地有

$$r \leq n. \quad (1.4.8)$$

定义 1.4.1 当 $r < n$, 即当矩阵 A 的线性无关的特征向量的个数小于矩阵的阶数时, 称 A 为 **亏损 (defective) 矩阵**;

反之，则称为非亏损的。

据此可知：可正规化的矩阵应是非亏损的。特别，当 A 的所有特征值均相异时，则 A 应是非亏损的。相似于正规阵、可正规化、非亏损这是不同说法的同一矩阵类，因其具有重要的应用，是需要着重进行研究的。

当 A 是非亏损矩阵时，由

$$H^{-1}AH = \text{diag}(\lambda_1, \dots, \lambda_n)$$

两端同用 D^{-1} 及 D 去左乘和右乘，这里 D 是非奇异的任意对角阵，可得

$$(HD)^{-1}A(HD) = \text{diag}(\lambda_1, \dots, \lambda_n), \quad (1.4.9)$$

适当选择 D 可使 HD 的各列成为单位长向量，它们都是 A 的特征向量，其个数为 n ，且彼此线性无关，此时称 A 是具有完全标准化特征向量系的矩阵。如果仍把 HD 记为 H ，因此，当 A 非亏损时，我们总可以假设，它具有变换矩阵 H ，其各列就是 A 的一组完全标准的特征向量组成的，记

$$H = (\mathbf{x}_1, \dots, \mathbf{x}_n), \|\mathbf{x}_k\| = 1, \quad (1 \leq k \leq n), \quad (1.4.10)$$

当 H 确定之后， H^{-1} 也随之而确定，记

$$(H^{-1})^T = (\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_n), \quad (1.4.11)$$

$$\frac{\mathbf{y}_k}{s_k} = \mathbf{g}_k, \quad |s_k| = 1/\|\mathbf{g}_k\|, \quad k = 1, 2, \dots, n,$$

则 $\mathbf{g}_k = (H^{-1})^T \mathbf{e}_k$, $\mathbf{g}_k^T = \mathbf{e}_k^T (H^{-1})$ ，由

$$H^{-1}A = \text{diag}(\lambda_1, \dots, \lambda_n)H^{-1}$$

可得

$$\mathbf{g}_k^T A = \lambda_k \mathbf{g}_k^T, \quad k = 1, 2, \dots, n$$

称 \mathbf{g}_k 为 A 的对应于特征值 λ_k 的左特征向量，而 $\mathbf{y}_k = \mathbf{g}_k/\|\mathbf{g}_k\|$ 为 A 的对应于特征值 λ_k 的标准化的左特征向量。由前面的分析可知，当 A 是非亏损的矩阵时，则存在 A 的右的和左的完全标准特征向量系。由于 $H^{-1}H = I$ ，故

$$e_i^T H^{-1} H e_j = g_i^T x_j = \frac{1}{s_i} y_i^T x_j = \begin{cases} 0, & i \neq j, \\ 1, & i = j \end{cases}$$

或

$$y_i^T x_j = \begin{cases} 0, & i \neq j, \\ s_i \neq 0, & i = j, \end{cases} \quad (1.4.12)$$

即非亏损阵的左右完全标准特征向量系形成一个双正交系。

在 (1.4.12) 中, 因 x_i, y_i 都是单位向量, 据 Cauchy 不等式, 有

$$|s_i| \leq 1, \quad i = 1, 2, \dots, n. \quad (1.4.13)$$

当 A 是正规阵时, A 的变换矩阵 H 可取为酉阵, 此时, 据 (1.4.11) 及 $H^{-1} = H^H$, 知有 $g_i = \bar{x}_i$ ($i = 1, \dots, n$), 从而,

$$|s_i| = |g_i| = 1 \quad (1.4.14)$$

即正规阵是具有完全标准化的双正交特征向量系的矩阵。

注意, 当 A 的特征值均相异时, 则 A 显然相似于以 A 的特征值为对角元的对角阵, 因而必为非亏损阵, 它的对应于某一特征值的特征向量, 如不计特征向量各分量间的比例因子, 应是唯一的, 其变换矩阵 H , 如不计列的次序, 则也是唯一的。但当 A 是特征值均不相异的非亏损阵时, 其变换矩阵 H 将不是唯一的, 即 A 可以有不同的变换矩阵 H 能同时使 A 成为对角阵。例如, 当 A 是单位阵 I 时, 实际上任意的非奇异矩阵都可作为 A 的变换矩阵。

反之, 如果同一的变换矩阵 H 可使两个不同的矩阵 A, B 同时化为对角阵, 自然地应该问 A, B 之间存在什么关系, 此时, 因

$$\begin{aligned} H^{-1} A H &= \text{diag}(\lambda_1, \dots, \lambda_n), \\ H^{-1} B H &= \text{diag}(\mu_1, \dots, \mu_n) \end{aligned} \quad (1.4.15)$$

将可推出 A, B 可换, 即

$$AB = BA, \quad (1.4.16)$$

进一步, 可得下列定理.

定理1.4.1 非亏损矩阵 A 、 B 可换的充要条件是, A 、 B 有相同的变换矩阵.

证明 只证条件的必要性.

由于 A 非亏损, 故存在非奇异阵 H , 使

$$H^{-1}AH = \text{diag}(\lambda_1, \dots, \lambda_n), \quad AH = H\text{diag}(\lambda_1, \dots, \lambda_n),$$

记

$$\text{diag}(\lambda_1, \dots, \lambda_n) = \text{diag}(\lambda_1 I_{r_1}, \lambda_2 I_{r_2}, \dots, \lambda_s I_{r_s}), \quad (1.4.17)$$

$$H = (\mathbf{h}_1^{(1)}, \dots, \mathbf{h}_{r_1}^{(1)}, \mathbf{h}_1^{(2)}, \dots, \mathbf{h}_{r_2}^{(2)}, \dots, \mathbf{h}_1^{(s)}, \dots, \mathbf{h}_{r_s}^{(s)}),$$

即设 A 中的相等的特征值都被排列成相邻的, (1.4.17) 中的诸 $\lambda_1, \dots, \lambda_s$ 视为相异的, 而 $\mathbf{h}_j^{(i)}$ 则表示 A 的对应于第 i 个相异特征值的第 j 个线性无关的特征向量. 由于

$$ABH = BH\text{diag}(\lambda_1 I_{r_1}, \dots, \lambda_s I_{r_s})$$

可见当 \mathbf{h} 是 A 的特征向量时, $B\mathbf{h}$ 也是 A 的特征向量, 可知:

$B\mathbf{h}_k^{(i)}$ 将属于 $\mathbf{h}_1^{(i)}, \dots, \mathbf{h}_{r_i}^{(i)}$ 张成的空间, 有

$$B\mathbf{h}_k^{(i)} = \sum_{j=1}^{r_i} \alpha_{jk}^{(i)} \mathbf{h}_j^{(i)}, \quad (1.4.18)$$

$$i = 1, \dots, s, \quad k = 1, \dots, r_i,$$

$$\therefore BH = H\text{diag}(X_1, \dots, X_s); \quad (1.4.19)$$

其中 $X_i(\alpha_{jk}^{(i)})$ 是 (1.4.18) 右端中的系数阵. 由于 B 非亏损, 故 $H^{-1}BH = \text{diag}(X_1, \dots, X_s)$ 应能化为对角阵, 即存在非奇异阵 $P_i (i = 1, \dots, s)$ 使

$$P_i X_i P_i^{-1} = \Lambda_i \quad (i = 1, \dots, s);$$

其中 Λ_i 为对角阵, 阶数为 r_i , ($i = 1, \dots, s$), 令

$$P = \text{diag}(P_1, \dots, P_s),$$

则 $(HP)^{-1}B(HP) = \text{diag}(\Lambda_1, \dots, \Lambda_s)$, 又

$$\begin{aligned} (HP)^{-1}AHP &= P^{-1}\text{diag}(\lambda_1 I_{r_1}, \dots, \lambda_s I_{r_s})P \\ &= \text{diag}(\lambda_1 I_{r_1}, \dots, \lambda_s I_{r_s}). \end{aligned}$$

可见 A 、 B 有相同的变换矩阵 HP 。

当 A 是亏损阵时，据 (1.4.7) 有定义：

定义 1.4.2 设 λ 是 A 的特征值， x 是非零向量， l 为某一正整数，如果

$$(A - \lambda I)^k x \neq 0, \quad k = 1, \dots, l-1;$$

但 $(A - \lambda I)^l x = 0$ ，则称 x 为 A 的对应于特征值 λ 的 l 级根向量。

据此，矩阵的特征向量便是矩阵的一级根向量，而矩阵的变换矩阵的各个列则是由矩阵的各级根向量组成的。由变换矩阵的非奇异性可知，这些各级根向量应是线性无关的。

由 (1.4.3) 可知，在矩阵的 Jordan 标准型中，同一个特征值可能对应着若干个 Jordan 块。设 $\lambda_1, \dots, \lambda_s$ 是 A 的相异特征值， r_i 是 λ_i 所对应的 Jordan 块中的阶数最高的 Jordan 块的阶，作多项式

$$\psi(\lambda) = (\lambda - \lambda_1)^{r_1} \cdots (\lambda - \lambda_s)^{r_s}, \quad (1.4.21)$$

可以证明

$$\psi(A) \equiv (A - \lambda_1 I)^{r_1} \cdots (A - \lambda_s I)^{r_s} = 0. \quad (1.4.22)$$

注意，若 A 的 Jordan 标准型中，对应某个特征值 λ_k ($1 \leq k \leq s$) 的 Jordan 块个数大于 1，此时对同一特征值的独立的特征向量个数将多于 1， $\psi(\lambda)$ 的次数必将小于矩阵 A 的阶数 n ，即必存在 A 的次数低于 n 的化零多项式 $\psi(\lambda)$ 。由此，我们给出下述定义。

定义 1.4.3 如果矩阵 A 对应某一特征值 λ 的线性无关的特征向量多于 1 时，则称矩阵是减次的，否则称为非减次的。

亏损与减次、非亏损与非减次是两个独立的概念，但矩阵 A 既非亏损又非减次的充要条件应是：矩阵 A 的诸特征值均相异。

§5 特征值扰动问题

5.1 矩阵特征值的连续相依性

设 A 为 n 阶方阵, 某特征值记为 λ , \tilde{A} 是 A 的近似矩阵, 其特征值记为 $\tilde{\lambda}$, 且设

$$\tilde{A} = A + \varepsilon B, \quad (1.5.1)$$

式中 ε 是实的或复的参数, B 是元素依模不大于 1 的矩阵。

所谓矩阵的特征值扰动问题, 就是要对 $|\lambda - \tilde{\lambda}|$ 的上界做出尽可能小的估计。在矩阵计算中, 由于初始输入的矩阵, 往往是某些物理量在仪器的量测精度范围内的测得的结果, 或是计算的经舍入处理后的结果, 它与真实的物理模型是有差别的。在 (1.5.1) 中, 如把 A 看作入机数据, \tilde{A} 看作真实的物理模型, 则 εB 便是它们之间的误差或扰动。可见讨论这种问题是具有普遍意义的。

由于 A 和 \tilde{A} 的特征值通常各有 n 个, 可见, $|\tilde{\lambda} - \lambda|$ 不仅与 ε 有关而且也 and A 、 \tilde{A} 的特征值的排序方法有关, 记 A 和 \tilde{A} 的特征多项式为

$$f(\lambda) = \det(\lambda I - A) = \lambda^n + a_1 \lambda^{n-1} + \cdots + a_n \quad (1.5.2)$$

$$\varphi(\varepsilon, \lambda) = \det(\lambda I - \tilde{A}) = \lambda^n + a_1(\varepsilon) \lambda^{n-1} + \cdots + a_n(\varepsilon) \quad (1.5.3)$$

这里 a_i ($i = 1, \cdots, n$) 是常数, $a_i(\varepsilon)$ ($i = 1, \cdots, n$) 是 ε 的次数不超过 i 的多项式, 满足

$$\lim_{\varepsilon \rightarrow 0} a_i(\varepsilon) = a_i \quad (i = 1, 2, \cdots, n) \quad (1.5.4)$$

定理 1.5.1 对任给正数 $\sigma > 0$, 存在 $\eta > 0$ 及 $f(\lambda) = 0$ 和 $\varphi(\varepsilon, \lambda) = 0$ 的根的排序方法, 使当 $|\varepsilon| < \eta$ 时, 有

$$|\tilde{\lambda}_k - \lambda_k| < (2n - 1)\sigma, \quad k = 1, \cdots, n \quad (1.5.5)$$

其中 λ_k 、 $\tilde{\lambda}_k$ 分别是 $f(\lambda) = 0$ 和 $\varphi(\varepsilon, \lambda) = 0$ 的第 k 个根。

证明 设 $\lambda_1, \dots, \lambda_n$ 是 $f(\lambda) = 0$ 的依任意次序排序的 n 个根, 当 λ_k 是它的 r 重根时, 则在 $j = 1, 2, \dots, n$ 中有 r 个下标 j_1, \dots, j_r 使 $\lambda_{j_s} = \lambda_k$ ($s = 1, \dots, r$); 对每个 λ_k , 以 λ_k 为中心以 σ 为半径作圆, 记为 C_k . 如果当 $k \neq j$ 时 C_k 与 C_j 有重叠部分, 则称 C_k 与 C_j , λ_k 与 λ_j 是相邻的; 如果对 $k \neq j$, 可以找到序列: $\lambda_{r_1}, \dots, \lambda_{r_s}$ 使

$$\lambda_k, \lambda_{r_1}, \dots, \lambda_{r_s}, \lambda_j$$

中的每一个数除 λ_k 外都与其前一个相邻, 称 λ_k 与 λ_j , C_k 与 C_j 是连通的. 于是诸圆 C_k ($k = 1, \dots, n$) 将在 λ 平面上组成有限个 (设为 m 个) 互不连通的连通域, 记为: G_1, G_2, \dots, G_m .

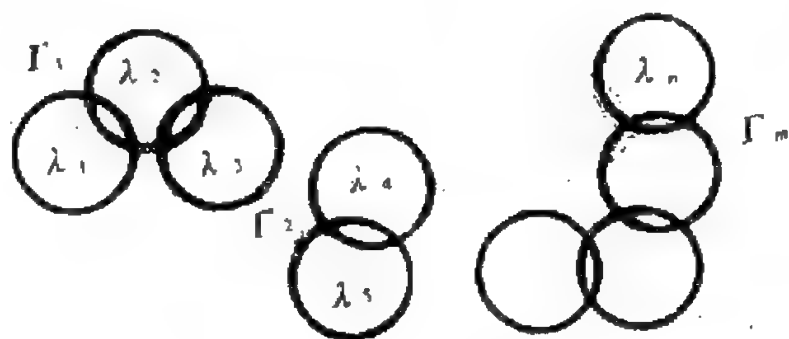


图1.1 σ -连通域

记 G_k 的边界为 Γ_k , 它们各由有限个圆弧构成, 令

$$\gamma = \max_{z \in \Gamma_k} |z|, \quad 1 \leq k \leq m$$

记
$$f_k = \min_{\lambda \in \Gamma_k} |f(\lambda)| > 0, \quad f = \min_{1 \leq k \leq m} (f_k)$$

由于

$$\varphi(\varepsilon, \lambda) = f(\lambda) + \sum_{k=0}^{n-1} (a_k(\varepsilon) - a_k) \lambda^k;$$

当 $\lambda \in \Gamma_k$ 时, 由 $|\lambda| \leq \gamma$, 有

$$|\varphi(\varepsilon, \lambda) - f(\lambda)| \leq \sum_{k=0}^{n-1} |a_k(\varepsilon) - a_k| |\lambda|^k.$$

据 $a_k(\varepsilon)$ 的连续性, 对任给 $\delta > 0$, 存在 $\eta_0 > 0$, 当 $|\varepsilon| \leq \eta_0$ 时,

有

$$|\varphi(\varepsilon, \lambda) - f(\lambda)| \leq \delta \cdot n \cdot \max(1, \gamma^{n-1});$$

当 $\delta < \delta_0 = f / (n \cdot \max(1, \gamma^{n-1}))$ 时, 知 $|\varphi(\varepsilon, \lambda) - f(\lambda)| < f \leq f_k$. 据 Rouché 定理, 此时, $\varphi(\varepsilon, \lambda)$, $f(\lambda)$ 在每个 Γ_k 围线内应有相同个数的零点; 这样, $f(\lambda)$, $\varphi(\varepsilon, \lambda)$ 的零点, 即 A , \tilde{A} 的特征值可按下列次序编排: 先依 G_k 的次序, 从 G_1 开始, 以任意次序分别对 $f(\lambda)$, $\varphi(\varepsilon, \lambda)$ 在围线内部的零点进行编号, 之后, 再对下一个 G_k 中的各自的零点进行编号, 为使符号简化, 重编以后的次序, 仍记为: $\lambda_1, \lambda_2, \dots, \lambda_n$ 和 $\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_n$. 任取 $1 \leq k \leq n$, 设 λ_k, λ_k 同属于某个 G_r , 且设 $f(\lambda)$ 在 G_r 中零点的个数为 $1 \leq m_r \leq n$, 因为 λ_k 是 G_r 中某个圆的圆心, 如果 $\tilde{\lambda}_k$ 所在圆的圆心是通过 s 个圆而与 λ_k 连通起来的, 于是

$$\begin{aligned} |\tilde{\lambda}_k - \tilde{\lambda}_k| &\leq |\lambda_k - \lambda_{\mu_1}| + |\lambda_{\mu_1} - \lambda_{\mu_2}| + \dots + |\lambda_{\mu_{s-1}} - \lambda_k| \\ &\leq \sigma + 2(s-1)\sigma \leq (2n-1)\sigma_\mu \end{aligned}$$

由于 σ 的任意性, 上述定理可以简单地表述为

$$\lim_{\varepsilon \rightarrow 0} \lambda_i(\varepsilon) = \lambda_i \quad (i = 1, 2, \dots, n) \quad (1.5.6)$$

或

$$|\tilde{\lambda}_i(\varepsilon) - \lambda_i| = o(|\varepsilon|^{p_i}) \quad (i = 1, \dots, n); \quad (1.5.7)$$

其中 $p_i > 0$, 而 $\lambda_1, \dots, \lambda_n, \lambda_1(\varepsilon), \dots, \lambda_n(\varepsilon)$ 各对应于 $A, A + \varepsilon B$ 的特征值和某一排序方法.

定义1.5.1 在 (1.5.7) 中, 当 $p_i \geq 1$ 时, 称矩阵 A 的特征值问题在 λ_i 处是**良态**的, 否则称为**病态**的.

例1.5.1 设

$$A = \begin{pmatrix} \lambda & 1 & & \\ & \lambda & 1 & \\ & & \ddots & \ddots \\ & & & 1 \\ & & & & \lambda \end{pmatrix},$$

矩阵 B 仅在 $(n, 1)$ 元素处为1, 其余均为零, 此时: 由

$$\det(\tilde{\lambda}I - A - \varepsilon B) = (\tilde{\lambda} - \lambda)^n + (-1)^n \varepsilon = 0,$$

得 $|\tilde{\lambda} - \lambda| = |\varepsilon|^{\frac{1}{n}},$

即 $p = \frac{1}{n} (< 1),$ A 的特征值问题是病态的.

5.2 Gerschgorin 定理

为了对矩阵特征值的性态问题做进一步的研究, 本段将以 Gerschgorin 圆盘定理及对角相似变换技术为工具来讨论这一问题.

设 n 阶方阵

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & & a_{nn} \end{pmatrix}, \quad (1.5.8)$$

对每个 $1 \leq i \leq n,$ 称

$$C_i = \{\lambda \mid |\lambda - a_{ii}| \leq \sum_{j \neq i} |a_{ij}| = \rho_i\} \quad (1.5.9)$$

为 A 的第 i 个 (Gerschgorin) 圆盘, 这里 a_{ij} 可为实数或复数.

定理 1.5.2 A 的任意一个特征值 λ 至少落在 A 的一个圆盘上.

证明 设 x 是 A 的对应于 λ 的特征向量, 且设其依模最大的分量为 $x_i = 1,$ 据 $Ax = \lambda x$ 的分量形式, 有

$$\sum_{j=1}^n a_{kj} x_j = \lambda x_k, \quad k = 1, 2, \cdots, n$$

当 $k = i$ 时, 得

$$\lambda - a_{ii} = \sum_{j \neq i} a_{ij} x_j,$$

故

$$|\lambda - a_{ii}| \leq \sum_{j \neq i} |a_{ij}|. \quad (1.5.10)$$

定理 1.5.3 设 A 的 n 个圆盘中, 有 s 个圆盘组成一个连

通区域 S ，且与其余的 $n-s$ 个圆盘互不相交，则 A 的 n 个特征值中，有且仅有 s 个落在 S 上。

证明 令 $A = D + C$ ，其中 D 是 A 的对角阵。作

$$A(\varepsilon) = D + \varepsilon C,$$

当 ε 由 0 增加到 1 时， $A(\varepsilon)$ 由 D 变到 A ，显然 $A(0)$ 的特征值就是 D 的对角元，它们分别是 n 个圆盘的圆心，其中有 s 个圆心在 S 上，由于 $A(\varepsilon)$ 的特征值连续依赖于 ε ，所以，当 ε 由 0 增加到 1 时， $A(\varepsilon)$ 的 n 个特征值将在 λ 平面上由各自的圆心出发画出 n 条连续曲线，这些曲线中的每一条或许全部落在 S 上，或许全部落在其余的 $n-s$ 个圆盘的交集上，由 S 上的 s 个圆的圆心出发的特征值，在 $0 < \varepsilon < 1$ 时应全部落在 S 上，而此种特征值有且仅有 s 个。

推论 1.5.1 设 A 有一个圆盘 C_i ，它孤立于其余的圆盘，则 $A(\varepsilon)$ 在此圆盘内有且仅有一个特征值 $\tilde{\lambda}(\varepsilon)$ ，当以 a_{ii} 作为 $\tilde{\lambda}(\varepsilon)$ 的近似值时，误差将不超过 $|\varepsilon| \sum_{j \neq i} |a_{ij}|$ ，即

$$|\tilde{\lambda}(\varepsilon) - a_{ii}| \leq \varepsilon \sum_{j \neq i} |a_{ij}|. \quad (1.5.11)$$

下面我们将利用上述定理和推论及对角相似变换技术来讨论矩阵特征值问题的性态。

不失一般性，设 λ_1 是任意矩阵 A 的一个特征值。先就 λ_1 是 A 的单重特征值的情形进行讨论。为方便起见，进一步再假设 A 是非亏损的矩阵。

我们的问题是要研究 A 的带扰动的矩阵 $\tilde{A} = A + \varepsilon B$ 的某一特征值 $\tilde{\lambda}_1(\varepsilon)$ 与 λ_1 的关系，这里假定当 $\varepsilon \rightarrow 0$ 时 $\tilde{\lambda}_1(\varepsilon) \rightarrow \lambda_1$ 。此时，由于存在非奇异矩阵 H 使

$$AH = H \text{diag}(\lambda_1, \dots, \lambda_n), \quad (1.5.12)$$

其中 $\lambda_1 \neq \lambda_i (i \geq 2)$ ，于是

$$H^{-1}(A + \varepsilon B)H = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) + \varepsilon H^{-1}BH$$

$$= \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix} + \varepsilon \begin{pmatrix} \beta_{11}/s_1 & \beta_{12}/s_1 & \dots & \beta_{1n}/s_1 \\ \beta_{21}/s_2 & \beta_{22}/s_2 & \dots & \beta_{2n}/s_2 \\ \vdots & \vdots & & \vdots \\ \beta_{n1}/s_n & \beta_{n2}/s_n & \dots & \beta_{nn}/s_n \end{pmatrix} \quad (1.5.13)$$

这里，我们假定 H 是由 A 的完全标准化特征向量系组成的，即

$$H = (\mathbf{x}_1, \dots, \mathbf{x}_n), \quad \|\mathbf{x}_i\| = 1 \quad (1 \leq i \leq n),$$

$$H^{-T} = (\mathbf{y}_1/s_1, \dots, \mathbf{y}_n/s_n), \quad \|\mathbf{y}_i\| = 1 \quad (1 \leq i \leq n), \quad (1.5.14)$$

$$s_i = \mathbf{y}_i^T \mathbf{x}_i \neq 0, \quad i = 1, \dots, n,$$

$$\mathbf{y}_i^T \mathbf{x}_j = 0, \quad i \neq j;$$

而

$$\beta_{ij} = \mathbf{y}_i^T B \mathbf{x}_j, \quad i, j = 1, \dots, n. \quad (1.5.15)$$

由于 $\lambda_1 \neq \lambda_i \quad (i = 2, \dots, n)$ ，故 (1.5.13) 右端的第 1 个圆盘，当 ε 充分小时，必将孤立于其余的圆盘，此时若以

$$\tilde{\lambda}_1'(\varepsilon) = \lambda_1 + \varepsilon \beta_{11}/s_1 \quad (1.5.16)$$

作为 $A + \varepsilon B$ 的一个近似特征值，则误差将不超过

$$\rho(\varepsilon) = |\varepsilon| \sum_{j=2}^n |\beta_{1j}/s_1| = O(|\varepsilon|),$$

但这一估计还可进一步精确化。为此，用 D^{-1} 和 D 分别左乘和右乘 (1.5.13) 的两端，其中取

$$D = \text{diag}(k/\varepsilon, 1, \dots, 1),$$

于是

$$D^{-1}H^{-1}(A + \varepsilon B)HD$$

$$= \text{diag}(\lambda_i) + \varepsilon \begin{pmatrix} \beta_{11}/s_1 & \varepsilon \beta_{12}/ks_1 & \dots & \varepsilon \beta_{1n}/ks_1 \\ k\beta_{21}/\varepsilon s_2 & \beta_{22}/s_2 & \dots & \beta_{2n}/s_2 \\ \vdots & \vdots & & \vdots \\ k\beta_{n1}/\varepsilon s_n & \beta_{n2}/s_n & \dots & \beta_{nn}/s_n \end{pmatrix} \quad (1.5.17)$$

此时，第 1 个和第 i 个圆盘的中心和半径为：

	中心	半径
第 1 个	$\lambda_1 + \varepsilon \beta_{11}/s_1$	$(\varepsilon^2/ ks_1) \sum_{j=2}^n \beta_{1j} $
第 i 个	$\tilde{\lambda}_i + \varepsilon \beta_{ii}/s_i$	$ k\beta_{ii}/s_i + \left \frac{\varepsilon}{s_i} \right \sum_{j=1, j \neq i}^n \beta_{ij} $

欲使第 1 个圆盘孤立于其余的圆盘，只需

$$\left| \lambda_1 - \lambda_i + \varepsilon \left(\frac{\beta_{11}}{s_1} - \frac{\beta_{ii}}{s_i} \right) \right| > |k\beta_{ii}/s_i| + \left| \frac{\varepsilon}{s_i} \right| \sum_{j=1, j \neq i}^n |\beta_{ij}| \\ + \left| \frac{\varepsilon^2}{ks_1} \right| \sum_{j=2}^n |\beta_{1j}|.$$

为此，可先取 k 使 $\max_{2 \leq i \leq n} |k\beta_{ii}/s_i| < \frac{1}{2} \min_{i \geq 2} |\lambda_1 - \lambda_i|$,

再取 ε 使

$$|\varepsilon| \left(\left| \frac{\beta_{11}}{s_1} \right| + \left| \frac{\beta_{ii}}{s_i} \right| + \left| \frac{1}{s_i} \right| \sum_{j=1, j \neq i}^n |\beta_{ij}| \right) + \left| \frac{\varepsilon}{ks_1} \right| \sum_{j=2}^n |\beta_{1j}| \\ < \frac{1}{2} \min_{i \geq 2} |\lambda_1 - \lambda_i|.$$

据 (1.5.17)，及 Gerschgorin 定理，则有

$$|\tilde{\lambda}_1(\varepsilon) - \tilde{\lambda}'_1(\varepsilon)| \leq \left| \frac{\varepsilon^2}{ks_1} \right| \sum_{j=2}^n |\beta_{1j}| = O(|\varepsilon|^2). \quad (1.5.18)$$

当 λ_1 是 A 的单重特征值但 A 本身是亏损矩阵时，在 A 的 Jordan 标准型中，对应于其余的特征值 $\lambda_i (i \geq 2)$ 将有阶数大于 1 的 Jordan 子块，此时，仿照上述的讨论过程，仍可得到 (1.5.18) 中的结果。

当 A 的 Jordan 标准型中，对应于 λ_1 的 Jordan 块中含有阶数大于 1 的块时，设相应的最高阶数为 r_1 ，此时，只能得到

$$|\tilde{\lambda}_1(\varepsilon) - \lambda_1| = O(|\varepsilon|^{\frac{1}{r_1}}) \quad (1.5.19)$$

的估计，即相应的特征值扰动问题是病态的。

§6 特征值问题的条件数

6.1 谱条件数

本节所讨论的矩阵 A 均指非亏损矩阵。

设 $\lambda_1, \lambda_2, \dots, \lambda_n$ 为 A 的特征值, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ 是相应的完全标准特征向量系。

$$A\mathbf{x}_i = \lambda_i \mathbf{x}_i \quad (i = 1, \dots, n), \quad (1.6.1)$$

记

$$H = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n), \quad (1.6.2)$$

令

$$H^{-T} = (\mathbf{y}_1/s_1, \dots, \mathbf{y}_n/s_n), \quad (1.6.3)$$

其中 $\|\mathbf{y}_i\| = 1$, 且

$$\mathbf{y}_j^T \mathbf{x}_i = \begin{cases} s_i \neq 0 & (i = j); \\ 0 & (i \neq j). \end{cases} \quad (1.6.4)$$

设 $\tilde{\lambda}(\epsilon)$ 是 $A + \epsilon B$ 的一个特征值, 则 $A + \epsilon B - \tilde{\lambda}I$ 是奇异阵, 知

$$H^{-1}(A + \epsilon B - \tilde{\lambda}I)H = \text{diag}(\lambda_i - \tilde{\lambda}) + \epsilon H^{-1}BH$$

亦应是奇异阵, 当 $\lambda_i \neq \tilde{\lambda} (i = 1, \dots, n)$ 时, 由于

$$\begin{aligned} & H^{-1}(A + \epsilon B - \tilde{\lambda}I)H \\ &= \text{diag}(\lambda_i - \tilde{\lambda}) [I - \epsilon (\text{diag}(\lambda_i - \tilde{\lambda}))^{-1} H^{-1}BH] \end{aligned}$$

可知

$$\|\epsilon (\text{diag}(\lambda_i - \tilde{\lambda}))^{-1} H^{-1}BH\| \geq 1 \quad (1.6.5)$$

对任意意义的矩阵范数成立, 而在欧氏范数意义下, 有

$$|\epsilon| \max_i (|\lambda_i - \tilde{\lambda}|^{-1}) \|H^{-1}\| \|H\| \|B\| \geq 1$$

或

$$\min_i |\lambda_i - \tilde{\lambda}| \leq \|\epsilon\| \|H^{-1}\| \|H\| \|B\|, \quad (1.6.6)$$

上式对 $\lambda_i = \tilde{\lambda}$ 也是成立的, 故在原讨论中关于 $\lambda_i \neq \tilde{\lambda}$ 的限制可以去掉, 记

$$\kappa(H) = \|H^{-1}\| \|H\|, \quad (1.6.7)$$

可知

$$|\lambda_i - \tilde{\lambda}| \leq |\varepsilon| \kappa(H) \|B\| \quad (1.6.8)$$

至少对 i 的一个值成立.

注意, 使得 $H^{-1}AH = \text{diag}(\lambda_1, \dots, \lambda_n)$ 的 H 一般是非唯一的, 如令

$$\kappa(A) = \inf_H (\kappa(H)), \quad (1.6.9)$$

显然, $\kappa(A)$ 能为某个 H 所取到, 且类似地有

$$|\lambda_i - \tilde{\lambda}| \leq |\varepsilon| \kappa(A) \|B\|. \quad (1.6.10)$$

由于 $\kappa(A)$ 是用矩阵的谱范数定义的, 且因 $\kappa(A)$ 的大小能刻画 A 的特征值扰动的敏感性, 因此, 通常称 $\kappa(A)$ 为 A 的对特征值问题来说的谱条件数.

当 $\kappa(A)$ 较小时, 人们有时也把 A 的特征值问题叫做是良态的, 而把相反的情形叫做是病态的. 但此地所谓的良态、病态与本章第 5 节定义的良好态与病态的概念实际上是有区别的. 由于本节只讨论非亏损阵的情形, 而非亏损阵的特征值问题在第 5 节中均应算作是良态的. 因此, 本节中提到的良态、病态概念实际是上一节中的良态概念的更细致的区分. 关于谱条件数 $\kappa(A)$, 可以证明它具有以下一些简单的性质:

$$(1) \kappa(A) \geq 1, \quad (1.6.11)$$

$$(2) \text{ 当 } A \text{ 是正规阵时, } \kappa(A) = 1. \quad (1.6.12)$$

由 (1.6.8) 还可建立类似于定理 1.5.3 的定理.

定理 1.6.1 如果圆盘

$$|\lambda - \lambda_i| \leq |\varepsilon| \kappa(A) \|B\|, \quad i = 1, 2, \dots, n \quad (1.6.13)$$

中有 s 个圆盘的和集组成与其余 $n-s$ 个圆盘互不连通的连通域

S , 则 $A + \varepsilon B$ 的特征值中恰有 s 个落在 S 上.

6.2 n -条件数

当 λ_1 是非亏损矩阵 A 的单重特征值时, 据 (1.5.16) 及 (1.5.17) 知道, 当以

$$\lambda_1'(\varepsilon) = \lambda_1 + \varepsilon \beta_{11}/s_1 \quad (1.6.14)$$

作为 $A + \varepsilon B$ 的一个特征值 $\lambda_1(\varepsilon)$ 的近似值时, 有

$$\lambda_1(\varepsilon) - \lambda_1 = \varepsilon \beta_{11}/s_1 + O(\varepsilon^2). \quad (1.6.15)$$

当 $|s_1|$ 很小时, 上式表明相应的特征扰动对于矩阵元素的扰动将是十分敏感的. 据此, 我们一般称

$$1/s_i, i = 1, 2, \dots, n \quad (1.6.16)$$

为矩阵 A 的 n -条件数. 诸 s_i 的计算可依 (1.4.10) ~ (1.4.12) 的公式进行.

谱条件数只能从总体上讨论特征扰动的敏感性而不能按特征值逐个地进行讨论, 因此, n -条件数的引入意在补救其不足之处.

6.3 n -条件数与谱条件数的关系

设 H 是 (1.6.9) 右端中取得准确下界的阵,

$$\mathbf{x}_i = H\mathbf{e}_i / \|H\mathbf{e}_i\|, \quad \mathbf{y}_i = H^{-T}\mathbf{e}_i / \|H^{-T}\mathbf{e}_i\|,$$

于是

$$\begin{aligned} |s_i| &= |\mathbf{y}_i^T \mathbf{x}_i| = \frac{|\mathbf{e}_i^T H^{-1} H \mathbf{e}_i|}{\|H\mathbf{e}_i\| \|H^{-T}\mathbf{e}_i\|} \\ &= \frac{1}{\|H\mathbf{e}_i\| \|H^{-T}\mathbf{e}_i\|} \geq \frac{1}{\kappa(A)}, \end{aligned}$$

即

$$\left| \frac{1}{s_i} \right| \leq \kappa(A), \quad i = 1, 2, \dots, n. \quad (1.6.17)$$

另一方面, 由

$$H^{-1}AH = \text{diag}(\lambda_1, \dots, \lambda_n),$$

$$H = (\mathbf{x}_1, \dots, \mathbf{x}_n), \quad H^{-T} = \left(\frac{\mathbf{y}_1}{s_1}, \dots, \frac{\mathbf{y}_n}{s_n} \right),$$

记

$$D = \text{diag}(|s_1|^{-\frac{1}{2}}, \dots, |s_n|^{-\frac{1}{2}}),$$

有

$$(HD)^{-1}A(HD) = \text{diag}(\lambda_1, \dots, \lambda_n),$$

据

$$\begin{aligned} \kappa(A) &\leq \|HD\| \|(HD)^{-1}\| \leq \|HD\|_E \|(HD)^{-1}\|_E \\ &\leq \left(\sum_{i=1}^n \frac{\|x_i\|^2}{|s_i|} \right)^{\frac{1}{2}} \left(\sum_{i=1}^n \frac{\|y_i\|^2}{|s_i|} \right)^{\frac{1}{2}} = \sum_{i=1}^n \frac{1}{|s_i|}, \quad (1.6.18) \end{aligned}$$

综合 (1.6.15) 和 (1.6.16), 有

$$\frac{1}{|s_i|} \leq \kappa(A) \leq \sum_{i=1}^n \frac{1}{|s_i|}. \quad (1.6.19)$$

这个不等式说明: 当 A 的某一特征值属于病态时, 则就总体上来说的特征值问题必是病态的, 而当总体的特征值问题属于病态时, 则在个别情况下的特征问题中至少有一个是病态的.

谱条件数和 n -条件数还有一个重要性质, 即二者在西相似变换下是不变的.

设 R 是酉阵, 且 $B = RAR^H$, 若 $H^{-1}AH = \text{diag}(\lambda_1, \dots, \lambda_n)$, 则有 $(RH)^{-1}B(RH) = \text{diag}(\lambda_1, \dots, \lambda_n)$, 记 B 的谱条件数 κ' , 则

$$\begin{aligned} \kappa' &= \inf_H (\|(RH)^{-1}\|_2 \|RH\|_2) \\ &= \inf_H (\|H^{-1}\|_2 \|H\|_2) = \kappa. \end{aligned} \quad (1.6.20)$$

若以 x_i, y_i 表示 A 的右、左特征向量, 则 B 的右、左特征向量为:

$$x'_i = Rx_i, \quad y'^T_i = e_i^T (RH)^{-1} / \|e_i (RH)^{-1}\|_2 = y_i^T R^H$$

于是

$$\begin{aligned} s'_i &= y'^T_i x'_i = y_i^T R^H R x_i = y_i^T x_i = s_i. \\ (i &= 1, 2, \dots, n) \end{aligned} \quad (1.6.21)$$

§7 实对称矩阵的扰动

7.1 实对称扰动

实对称阵、Hermite 阵及正规阵之间存在着前者包含于后者的关系，且均属于非亏损的矩阵，它们的特征扰动问题必有更深入的结果。本节讨论实对称阵的性质及其特征扰动问题，有关的结论常可推之于后面的两类矩阵，但在叙述时将不一一指出。

设 A 是实对称阵，其特征值为 $\lambda_i (i=1, 2, \dots, n)$ ，它们对应的右的和左的标准化特征向量设为 x_i 和 y_i ，不论诸 λ_i 是否相异，可取 $x_i = y_i (i=1, \dots, n)$ 为实分量的向量，于是有

$$s_i = y_i^T x_i = 1 \quad (i=1, \dots, n), \quad (1.7.1)$$

因此，实对称阵的特征值扰动问题总是良态的。记 $\lambda(\varepsilon)$ 为 $A + \varepsilon B$ 的任意一个特征值，仿本章 §6(1.6.1) ~ (1.6.8) 的证明，可知

$$|\lambda(\varepsilon) - \lambda_i| \leq |\varepsilon| \|B\| \leq |\varepsilon| \|B\|_F \leq n|\varepsilon| \quad (1.7.2)$$

至少对某一 $i=1, 2, \dots, n$ 成立。即 $A + \varepsilon B$ 的任意一个特征值必将落在某一个以 λ_i 为中心以 $n|\varepsilon|$ 为半径的圆盘上。由此，可以导出一个有趣的结论：如果 λ_i 是 A 的单重实特征值，则当 ε 充分小时，只要 εB 是实的但不一定对称， $A + \varepsilon B$ 必有一个靠近 λ_i 的实特征值。

因为 λ_i 既是实数又是单重特征值，于是可以假设 A 有一个异于 λ_i 又最靠近 λ_i 的特征值 λ_k ，当 $|\varepsilon| < |\lambda_i - \lambda_k|/2n$ 时，设 $\lambda(\varepsilon)$ 是 $A + \varepsilon B$ 的且落在 $|\lambda - \lambda_i| < n|\varepsilon|$ 圆盘上的特征值，当 $k \neq i$ 时，由

$$|\lambda_k - \lambda(\varepsilon)| \geq |\lambda_k - \lambda_i| - |\lambda_i - \lambda(\varepsilon)| > n|\varepsilon|. \quad (1.7.3)$$

此外，如果 A 的特征值均相异，则当 $|\varepsilon|$ 充分小且 εB 为实阵时， $A + \varepsilon B$ 的特征值亦将全是实数且相异，此时 $A + \varepsilon B$ 应有完全的特征向量系。当 B 不仅是实的且也是对称阵时，不论实参数 ε 的大小如何，扰动后的矩阵仍是实对称矩阵，在下面的讨论中，我们将把 ε 取为 1，并把 $|B| \leq 1$ 的限制去掉，而提出更一般的问题。

设 A, B 均为实对称阵, 记

下面分两种情况进行讨论.

$$C = \begin{pmatrix} \alpha & \vdots & \alpha^T \\ \dots\dots\dots & & \\ \alpha & \vdots & \text{diag}(\alpha) \end{pmatrix}. \quad (1.7.5)$$
$$x_1 > y_1 > x_2 > y_2 > \dots x_n > y_n > x_{n+1} \quad (1.7.6)$$

定理1.7.1 当 (1.7.5) 中的 a_i 均相异且 a 的分量均不为零时, 则 C 的特征值: $\lambda_1 > \lambda_2 > \dots > \lambda_n$ 与诸 a_i 在严格的意义下是隔离的.

• 31 •

$> \cdots > a_{n-1}$, 此时 C 的特征方程为

$$(a - \lambda) \prod_{i=1}^{n-1} (a_i - \lambda) - \sum_{i=1}^{n-1} a_i^2 \prod_{\substack{j=1 \\ j \neq i}}^{n-1} (a_j - \lambda) = 0, \quad (1.7.7)$$

记

$$f(\lambda) = \lambda + \sum_{i=1}^{n-1} a_i^2 / (a_i - \lambda), \quad (1.7.8)$$

则 (1.7.7) 化为

$$a - \lambda - \sum_{i=1}^{n-1} a_i^2 / (a_i - \lambda) = a - f(\lambda) = 0. \quad (1.7.9)$$

如图1.2, 因

$$f'(\lambda) = 1 + \sum_{i=1}^{n-1} a_i^2 / (a_i - \lambda)^2 > 0, \quad (1.7.10)$$

$$a_{i+1} < \lambda < a_i, \quad i = 1, 2, \dots, n-2.$$

方程 (1.7.9) 在 (a_{i+1}, a_i) 内有根, 记为 λ_{i+1} , 在 $(-\infty, a_n)$ 和 (a_1, ∞) 的根记为 λ_n 和 λ_1 则 $\lambda_1, \lambda_2, \dots, \lambda_n$ 与 a_1, \dots, a_{n-1} , 在严格的意义下是隔离的, 即 $\lambda_1 > a_1 > \lambda_2 > \cdots > a_{n-1} > \lambda_n$.

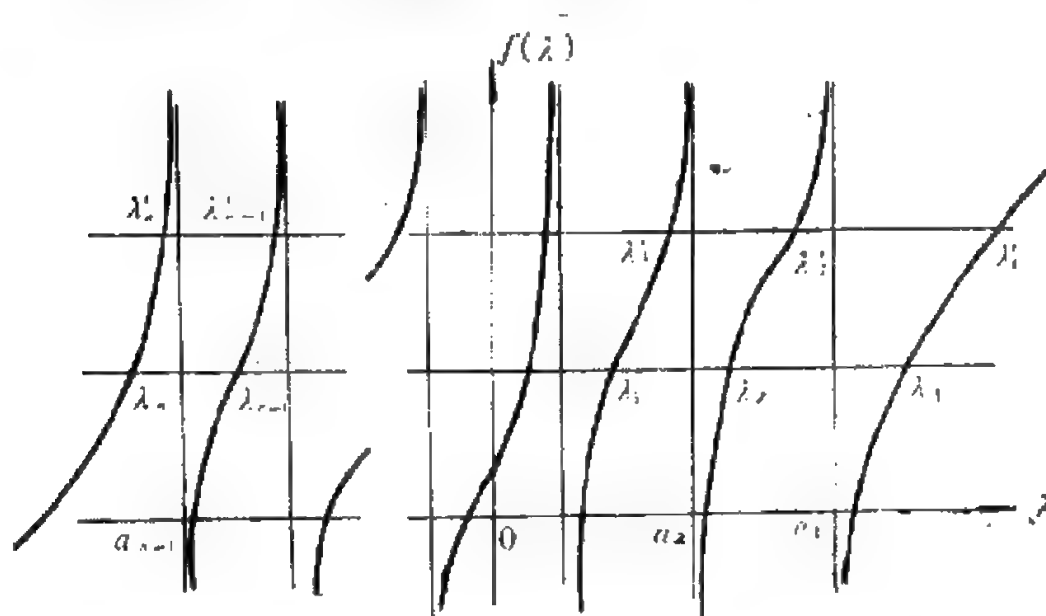


图1.2 特征值的隔离

当某个 a_i 是 $\text{diag}(a_1, \dots, a_n)$ 的 r_i 重特征值时, 由

$$a_{i-1} > a_i = a_{i+1} = \cdots = a_{i+r_i-1} > a_{i+r_i}$$

可见

$$\alpha_i = \lambda_{i+1} = \alpha_{i+1} = \cdots = \lambda_{i+r_i-1} = \alpha_{r_i+i-1},$$

即 λ_{i+1} 必为 C 的 r_{i-1} 重特征值, 此时, $\lambda_1, \cdots, \lambda_n$ 与 $\alpha_1, \cdots, \alpha_{n-1}$ 在弱的意义下是隔离的。

当 α 的某些分量是零时, 则与之同行的 α_i 必为 C 的特征值, 这些 α_i 将既属于 C 的特征值数组, 又属数组 $\alpha_1, \cdots, \alpha_{n-1}$, 把 C 中这些 α_i 所在的行列去掉, 据前面的证明知, 剩下的矩阵的特征值将在弱的意义下为 $\text{diag}(\alpha_i)$ 的剩下的那些 α_i 所隔离。而在弱的意义下互相隔离的两个数组中, 各自分别插入一些相同的数后, 仍将在弱的意义下相互隔离。于是, 不管诸 α_i 是否相异, 也不管 α 是否有零分量, C 的特征值将在弱的意义下为 $\alpha_1, \cdots, \alpha_{n-1}$ 所隔离。

当 α 及诸 α_i 均不变, 而 a 从 α 变到 α' 时, 设 C 的特征值由 λ_i 变到 λ'_i , 而 λ_i 与 α_i 的隔离情况与 λ'_i 和 α_i 的隔离情况是一样的。因此, λ_i 中的重特征值与 λ'_i 中的重特征值将不变化。当 λ_i, λ'_i 都是单重特征值时, 由于

$$\begin{aligned} f(\lambda'_k) - f(\lambda_k) &= \lambda'_k + \sum_{i=1}^{n-1} \alpha_i^2 / (\alpha_i - \lambda'_i) - \lambda_k - \sum_{i=1}^{n-1} \alpha_i^2 / (\alpha_i - \lambda_i) \\ &= \alpha' - \alpha, \end{aligned} \quad (1.7.11)$$

有

$$(\lambda'_k - \lambda_k) f'(\xi_k) = \alpha' - \alpha, \quad \lambda'_k - \lambda_k = m_k (\alpha' - \alpha), \quad (1.7.12)$$

这里 ξ_k 在 λ_k 与 λ'_k 之间, $0 < m_k = 1/f'(\xi_k) < 1$ 。由 (1.7.7) 及根与系数的关系, 得

$$\sum \lambda_k = \alpha + \sum \alpha_i, \quad \sum \lambda'_k = \alpha' + \sum \alpha_i,$$

故

$$\sum \lambda'_k - \sum \lambda_k = \alpha' - \alpha,$$

从而

$$\sum m_k = 1. \quad (1.7.13)$$

由 (1.7.12) 和 (1.7.13) 可知对角加边矩阵 C 中, 当 α 发生 $\Delta\alpha = \alpha' - \alpha$ 的改变量时, 特征值移动的方向与 $\Delta\alpha$ 的方向相同, 且特征值的总变化量等于 $\Delta\alpha$.

7.2.2 $C = A + B$, B 是秩 1 矩阵.

因为 B 是秩 1 的, 存在正交阵 R 使

$$B = R \begin{pmatrix} \rho & \vdots & 0 \\ \cdots & \cdots & \cdots \\ 0 & \vdots & 0 \end{pmatrix} R^T. \quad (1.7.14)$$

在 $C = A + B$ 的两端分别用 R 、 R^T 左乘和右乘, 记

$$RAR^T = \begin{pmatrix} \alpha & \vdots & \mathbf{a}^T \\ \cdots & \cdots & \cdots \\ \mathbf{a} & \vdots & A_{n-1} \end{pmatrix}.$$

这里 A_{n-1} 是 $n-1$ 阶的实对称阵, 得

$$RCR^T = \begin{pmatrix} \alpha + \rho & \vdots & \mathbf{a}^T \\ \cdots & \cdots & \cdots \\ \mathbf{a} & \vdots & A_{n-1} \end{pmatrix}.$$

对 A_{n-1} , 存在 $n-1$ 阶实正交阵 S , 使

$$SA_{n-1}S^T = \text{diag}(\alpha_1, \dots, \alpha_{n-1}),$$

作

$$Q = \begin{pmatrix} 1 & \vdots & 0 \\ \cdots & \cdots & \cdots \\ 0 & \vdots & S \end{pmatrix} R$$

则

$$Q(A+B)Q^T = \begin{pmatrix} \alpha + \rho & \vdots & \mathbf{b}^T \\ \cdots & \cdots & \cdots \\ \mathbf{b} & \vdots & \text{diag}(\alpha_i) \end{pmatrix},$$

这里 $\mathbf{b} = S\mathbf{a}$, 上式右端已化为对角加边的实对称阵, 由于 A 和 $A+B$ 的特征值就是

$$\begin{pmatrix} \alpha & \vdots & \mathbf{b}^T \\ \cdots & \cdots & \cdots \\ \mathbf{b} & \vdots & \text{diag}(\alpha_i) \end{pmatrix} \text{ 和 } \begin{pmatrix} \alpha + \rho & \vdots & \mathbf{b}^T \\ \cdots & \cdots & \cdots \\ \mathbf{b} & \vdots & \text{diag}(\alpha_i) \end{pmatrix}$$

的特征值, 设其各为 λ_i 和 λ'_i , 据 (1.7.12), 有

$$\begin{aligned}\lambda'_i - \lambda_i &= m_i \rho, \\ \sum m_i &= 1 \quad (0 \leq m_i \leq 1).\end{aligned}\tag{1.7.15}$$

注意, 由于 ρ 是 B 的唯一的非零特征值, 上式表明: 当对称阵的扰动阵是秩 1 的对称阵时, 则特征值的扰动与扰动阵的唯一的非零特征值同号, 且特征值扰动之和等于扰动阵的特征值。

§8 Hermite 阵的极值性质

8.1 Hermite 阵的极值性质

本节设 A 为 Hermite 阵, 其特征值是实数:

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n.\tag{1.8.1}$$

存在酉阵

$$U = (\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n)\tag{1.8.2}$$

使

$$U^H A U = \text{diag}(\lambda_1, \cdots, \lambda_n).\tag{1.8.3}$$

视 \mathbf{x}_i ($i = 1, \cdots, n$) 为 n 维复向量空间 \mathbf{C}^n 的一组标准正交基。

对任意的 $\mathbf{x} \in \mathbf{C}^n$, \mathbf{x} 可以分别表为

$$\begin{aligned}\mathbf{x} &= \xi_1 \mathbf{e}_1 + \xi_2 \mathbf{e}_2 + \cdots + \xi_n \mathbf{e}_n; \\ \mathbf{x} &= \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \cdots + \alpha_n \mathbf{x}_n.\end{aligned}\tag{1.8.4}$$

当 $\mathbf{x}^H \mathbf{x} = 1$ 时, 可得

$$\begin{aligned}\mathbf{x}^H \mathbf{x} &= |\xi_1|^2 + |\xi_2|^2 + \cdots + |\xi_n|^2 \\ &= |\alpha_1|^2 + |\alpha_2|^2 + \cdots + |\alpha_n|^2 = 1.\end{aligned}\tag{1.8.5}$$

由于

$$\begin{aligned}f(\mathbf{x}) &= \mathbf{x}^H A \mathbf{x} = \sum_{i,j=1}^n a_{ij} \bar{\xi}_i \xi_j, \\ a_{ij} &= \bar{a}_{ji} \quad (i, j = 1, \cdots, n),\end{aligned}\tag{1.8.6}$$

易知 $f(\mathbf{x})$ 是 \mathbb{C}^n 上的实的连续泛函，它在球面

$$\mathbf{x}^H \mathbf{x} = 1 \quad (1.8.7)$$

取得最大值

$$\rho_1 = \max_{\|\mathbf{x}\|=1} (\mathbf{x}^H \mathbf{A} \mathbf{x})$$

由于

$$\mathbf{x}^H \mathbf{A} \mathbf{x} = \sum_{i=1}^n \lambda_i |\alpha_i|^2 \leq \lambda_1 \|\mathbf{x}\|^2 = \lambda_1,$$

故

$$\rho_1 \leq \lambda_1.$$

另外，由

$$\mathbf{x}_1^H \mathbf{A} \mathbf{x}_1 = \lambda_1 \|\mathbf{x}_1\|^2 = \lambda_1,$$

又有

$$\rho_1 \geq \lambda_1,$$

得

$$\rho_1 = \max_{\|\mathbf{x}\|=1} (\mathbf{x}^H \mathbf{A} \mathbf{x}) = \lambda_1. \quad (1.8.8)$$

当对 \mathbf{x} 不仅给以 $\mathbf{x}^H \mathbf{x} = 1$ 的约束，而且给以 $\mathbf{x}_1^H \mathbf{x} = 0$ 的约束时，可证

$$\max_{\substack{\|\mathbf{x}\|=1 \\ \mathbf{x}_1^H \mathbf{x} = 0}} (\mathbf{x}^H \mathbf{A} \mathbf{x}) = \lambda_2. \quad (1.8.9)$$

一般地，可证：

$$\max_{\substack{\|\mathbf{x}\|=1 \\ \mathbf{x}_i^H \mathbf{x} = 0 \\ i=1, \dots, k}} (\mathbf{x}^H \mathbf{A} \mathbf{x}) = \lambda_{k+1} \quad (1 \leq k+1 \leq n). \quad (1.8.10)$$

在上述讨论中，当要确定 \mathbf{A} 的第 $i+1$ 个特征值时，应先知道对应于 \mathbf{A} 的前 i 个特征值 $\lambda_1, \dots, \lambda_i$ 的那些特征向量 $\mathbf{x}_1, \dots, \mathbf{x}_i$ ，而利用下面的讨论则可避免这一点。

今考虑求 $\mathbf{x}^H \mathbf{A} \mathbf{x}$ 的极小极大极小问题：

$$\min_{\substack{P_i, \|x\|=1 \\ P_i^H x=0, P_i \neq 0 \\ i=1, \dots, k}} \max (x^H A x), \quad (1.8.11)$$

这里 P_1, \dots, P_k 为任意的非零向量. 仍令 $U = (x_1, \dots, x_n)$, 且记 $x = Uy$, 则 (1.8.11) 化为

$$\min_{\substack{q_i, \|y\|=1 \\ q_i^H y=0, q_i \neq 0 \\ i=1, \dots, k}} \max y^H \text{diag}(\lambda_i) y. \quad (1.8.12)$$

这里

$$q_i = U^H P_i, \quad i = 1, 2, \dots, k, \quad (1.8.13)$$

因为对于问题 (1.8.12), 存在一组向量 $\bar{q}_i (i = 1, \dots, k)$ 使

$$\begin{aligned} \min_{\substack{q_i, \|y\|=1 \\ q_i^H y=0, q_i \neq 0 \\ i=1, \dots, k}} \max y^H \text{diag}(\lambda_i) y &= \max_{\substack{\|y\|=1 \\ \bar{q}_i^H y=0, \bar{q}_i \neq 0 \\ i=1, \dots, k}} y^H \text{diag}(\lambda_i) y \\ &\geq \max_{\substack{\|y\|=1, \quad \eta_1=\dots=\eta_k=0 \\ \bar{q}_i^H y=0, \quad q_i \neq 0 \\ i=1, \dots, k}} y^H \text{diag}(\lambda_i) y = \max_{\substack{\|y\|=1 \\ \bar{q}_i^H y=0, q_i \neq 0 \\ i=1, \dots, k}} \sum_{i=k+1}^n \lambda_i |\eta_i|^2 \geq \lambda_{k+1}. \end{aligned}$$

这里 η_i 是 y 的第 i 个分量. 另一方面,

$$\begin{aligned} \min_{\substack{q_i, \|y\|=1 \\ q_i^H y=0, q_i \neq 0 \\ i=1, \dots, k}} \max y^H \text{diag}(\lambda_i) y &\leq \max_{\substack{\|y\|=1 \\ q_i^H y=0 \\ i=1, \dots, k}} y^H \text{diag}(\lambda_i) y \\ &= \sum_{i=k+1}^n \lambda_i |\eta_i|^2 \leq \lambda_{k+1}, \end{aligned}$$

最后得

$$\min_{\substack{q_i, \|y\|=1 \\ q_i^H y=0, q_i \neq 0 \\ i=1, \dots, k}} \max y^H \text{diag}(\lambda_i) y = \lambda_{k+1},$$

即

$$\min_{\substack{P_i, \|x\|=1 \\ P_i^H x=0, P_i \neq 0 \\ i=1, \dots, k}} \max x^H A x = \lambda_{k+1}. \quad (1.8.14)$$

8.2 Hermite 阵的和的特征值

设 A 、 B 、 C 均为 Hermite 阵, α_i 、 β_i 、 γ_i 分别是它们的按非增次序排列的特征值, 设

$$C = A + B,$$

由于

$$\gamma_k = \min_{\substack{P_i \\ |x|=1 \\ P_i^H x=0, P_i \neq 0 \\ i=1, \dots, k-1}} \max (x^H C x), \quad (1.8.15)$$

仍设 $U^H A U = \text{diag}(\alpha_i)$, 且取 P_i 属于 U 的前 $k-1$ 列, 即取 $P_i = U e_i$. 记 $x = U y$, 由 (1.8.15)

$$0 = P_i^H x = e_i^T U^H U y = e_i^T y,$$

$$x^H A x = y^H \text{diag}(\alpha_i) y = \sum_{i=1}^n \alpha_i |\eta_i|^2 \leq \alpha_k,$$

故

$$\gamma_k \leq \max_{\substack{|x|=1 \\ P_i^H x=0, P_i \neq 0 \\ i=1, \dots, k-1}} (x^H C x) \leq \alpha_k + \beta_1. \quad (1.8.16)$$

由于 $A = C + (-B)$, 视 C 、 $-B$ 、 A 为原来的 A 、 B 、 C , 类比 (1.8.16) 可得

$$\alpha_k \leq \gamma_k - \beta_n, \quad (1.8.17)$$

故

$$\alpha_k - \beta_n \leq \gamma_k \leq \alpha_k + \beta_1. \quad (1.8.18)$$

当 B 的元素满足 $|b_{ii}| \leq \varepsilon$ 时, 则由

$$-n\varepsilon \leq \beta_n < \beta_1 \leq n\varepsilon$$

可得

$$|\gamma_i - \alpha_i| \leq n\varepsilon, \quad i = 1, \dots, n. \quad (1.8.19)$$

利用极大极小原理还可进一步证明

$$\gamma_{r+s-1} \leq \alpha_r + \beta_s, \quad \gamma + s - 1 \leq n. \quad (1.8.20)$$

下面我们还将介绍一个有关 n 阶实对称矩阵和它的 $n-1$ 阶主子阵的特征值的定理.

定理1.8.1 设 A 为 n 阶实对称阵, A_{n-1} 是它的 $n-1$ 阶首主子阵, 若以 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$, $\lambda'_1 \geq \lambda'_2 \geq \dots \geq \lambda'_{n-1}$ 分别表示它们的特征值, 则诸 λ_i 为诸 λ'_i 在弱的意义下所隔离.

证明 记 $I' = \text{diag}(1, \dots, 1, 0)$,

$$x^T = (\xi_1, \dots, \xi_n) = x'^T + (0, \dots, 0, \xi_n),$$

$$A = A' + \begin{pmatrix} 0 & \vdots & a \\ \dots\dots\dots & & \\ a^T & \vdots & a_{nn} \end{pmatrix}.$$

考察

$$\min_{\substack{P_i: \|x'\|=1 \\ P_i^H x' = 0, P_i \neq 0 \\ i=1, \dots, s-1}} \max x'^H A x' = \min_{\substack{P_i: \|x'\|=1 \\ P_i^H x' = 0, P_i \neq 0 \\ i=1, \dots, s-1}} \max x'^H A' x' = \lambda'_s,$$

(1.8.21)

而左端又可写为

$$\min_{\substack{P_i: \|x\|=1 \\ P_i^H x = 0, P_i \neq 0 \\ i=1, \dots, s-1 \\ e_n^T x = 0}} \max x^H A x \geq \min_{\substack{P_i: \|x\|=1 \\ P_i^H x = 0, P_i \neq 0 \\ i=1, \dots, s}} \max (x^H A x) = \lambda_{s+1},$$

故

$$\lambda'_s \geq \lambda_{s+1}.$$

在 (1.8.21) 右端中, 当换 x' 为 x 时, 其值应是非减的, 即

$$\lambda_s \geq \lambda'_s.$$

最后, 有

$$\lambda_{s+1} \leq \lambda'_s \leq \lambda_s, \quad (1.8.22)$$

即 λ_i 被 λ'_i 在弱的意义下所隔离.

第一章 习 题

1.1 试将矩阵

$$\begin{pmatrix} 1 & 2 & 3 \\ 2 & 3 & 5 \\ 3 & 4 & 7 \end{pmatrix}, \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}$$

分解成等秩的列矩阵和行矩阵的乘积.

1.2 设 $A \in \mathbb{C}^{n \times m}$ 且 A 是秩 1 的, 试证存在非零的向量 $u \in \mathbb{C}^n$, $v \in \mathbb{C}^m$ 使 $A = uv^T$, 此种 u 、 v 唯一否?

1.3 求矩阵

$$A = \begin{pmatrix} 7 & 1 & & \\ & 7 & 0 & \\ & & 7 & 1 \\ & & & 7 \end{pmatrix}$$

的线性独立的各级根向量. 矩阵 A 是亏损阵吗? 是减次的吗?

1.4 设

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -\alpha_1 & -\alpha_2 & -\alpha_3 & -\alpha_4 \end{pmatrix},$$

问 A 可否有形如 $x^T = (0, x, x, x)$ 的特征向量? 如其有形如 $x^T = (1, x, x, x)$ 的特征向量, 试求之. 给出 A 非亏损的一个充分条件.

1.5 通过对行列式的直接展开, 证明 $E = I - \delta uv^H$ 的 $\det(E) = 1 - \delta v^H u$.

1.6 用 Gerschgorin 定理说明矩阵

$$\begin{pmatrix} 9 & 1 & -2 & 1 \\ 0 & 8 & 1 & 1 \\ -1 & 0 & 7 & 0 \\ 1 & 0 & 0 & 1 \end{pmatrix}$$

至少有二实根.

1.7 设 $B \in \mathbb{C}^{n \times n}$ 适合 $|b_{ii}| > \sum_{j=1, j \neq i}^n |b_{ij}|, (i=1, \dots, n)$,

试证 $\det(B) \neq 0$.

1.8 根据特征值的隔离定理说明

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 5 \\ 3 & 5 & 6 \end{pmatrix}$$

的特征值的分布情况.

1.9 设 H 是实对称阵, 它在正交相似变换下化为

$$R^T H R = \begin{pmatrix} 1.02 & -0.3 & 0.1 \\ -0.3 & 0.87 & -0.1 \\ 0.1 & -0.1 & 7.23 \end{pmatrix},$$

试估计 H 的特征值.

1.10 就矩阵

$$H = \begin{bmatrix} 4 & 3i \\ -3i & 2 \end{bmatrix}$$

验证极值原理.

1.11 设 $A \in \mathbb{C}^{n \times n}$ 且 A 的各列是互相正交的单位向量, 证明 $\max_{\|u\|=1} |(Au, u)|$ 只能在 u 是对应于 A 的依模最大的特征值的特征向量时取得.

1.12 证明任意方阵可经由酉相似变换化为上三角阵.

1.13 证明同时是酉阵又是三角阵的矩阵必是对角阵, 且对角元依模等于1, 证明同时是正规阵又是三角阵的矩阵必定是对角阵.

1.14 设 $F \in \mathbb{C}^{n \times m} (n \geq m)$, F 的秩为 m , 证明存在唯一的 $P_F \in \mathbb{C}^{n \times n}$, 使对任意的 $x \in \mathbb{R}^n$ 有

$$x = P_F x + (x - P_F x),$$

其中 P_F 是 F 的列的组合, $(x - P_F x)$ 与 F 的任意的列正交, 试导出 P_F 的表达式.

1.15 设 $T = (t_{ij})$ 是上三角阵, 当 $t_{ii} \neq t_{jj}$ 时, 可以适当选择 δ 使当 $i > j$ 时, 有

$$e_i E(e_i, e_j - \sigma) T E(e_i, e_j, \sigma) = 0,$$

从而证明对任意矩阵 A , 存在矩阵 $X = \Omega$, 其中 ΩF 是初等 Hermite 阵的乘积, F 是形如 $E(e_i, e_j, \sigma)$ 的初等矩阵的乘积, 使得

$$X^{-1} A X = \text{diag}(R_1, R_2, \dots),$$

其中每个 R_i 是三角阵且其对角元是彼此相等的.

1.16 幂零矩阵只能有零特征值, 幂等矩阵只能有 0 和 1 的特征值, 此处幂零和幂等矩阵分别是指满足方程 $A^v = 0 \neq A^{v-1}$ 和 $A^2 = A$ 的矩阵, 这里 v 是某一正整数, 称为 A 的幂零指数.

1.17 设 A 的幂零指数是 v , 则存在向量 v_1 , 其 Krylov 序列 $v_1, v_2, \dots, v_v \dots$ 定义为

$$v_{k+1} = A v_k \quad (k = 1, 2, \dots)$$

满足 $v_{i+1} \neq 0$, 记 $V = (v_1, \dots, v_v)$ 则有

$$A V = V J,$$

1.18 设 A, B 均是可正规化矩阵则下之二关系等价:

(1) $AB = BA$,

(2) 存在非奇异阵 X 使 $X^{-1} A X, X^{-1} B X$ 同是对角阵.

1.19 用两种方法展开 $\det \begin{pmatrix} \lambda I & -A \\ -B & \lambda I \end{pmatrix}$ 证明 AB, BA 的特征多项式除一个 λ 幂因子外是相等的.

1.20 设 A, B 均为方阵, 证明存在 $X \neq 0$ 满足 $AX = XB$ 的充要条件是 A, B 至少有一个公共的特征值, 据此证明 $AX - XB = C \neq 0$ 有唯一解的充要条件是 A, B 没有公共特征值.

1.21 如果 A^v 的迹数对一切正整数 $v > 0$ 均为零, 则 A 必为幂零矩阵, 反之也成立.

1.22 设 $A = J + \varepsilon^n e_1 e_n^T$, 其中

$$J = \begin{pmatrix} 0 & & & \\ 1 & 0 & & \\ & 1 & \ddots & \\ & & \ddots & 0 \\ & & & 1 & 0 \end{pmatrix},$$

证明 A 的特征值为 $\omega^k \varepsilon$, 这里 $\omega = \exp(i\pi/n)$.

1.23 设 $F = I - f e e^T$, $f^T = (f_1, \dots, f_n)$, 验证

$$\det(F - \lambda I) = \begin{vmatrix} -\lambda & -\lambda^2 & \cdots & \lambda^{n-1} & -\varphi(\lambda) \\ 1 & 0 & & & \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & \\ & & & & 0 \end{vmatrix},$$

这里 $\varphi(\lambda) = \lambda^n + f_n \lambda^{n-1} + \cdots + \lambda f_2 + f_1$.

1.24 与对角阵可换的矩阵的一般形式是什么与 Jordan 块可换的矩阵的一般形式又是什么?

1.25 任意实矩阵 A 可唯一地分解成一个实对称阵 B 和一个实反对称阵 C 之和, 而 A 为实正规阵的充要条件是 B, C 可换.

参 考 书

1. J. H. Wilkinson "The Algebraic Eigenvalue Problems". England, Oxford Univ. Press, 1965.
2. Alston S. Householder "The Theory of Matrices in Numerical Analysis". Blaisdell, New York, 1964.
3. Joel N. Franklin. "Matrix Theory". Prentice, Inc., Englewood Cliffs, New Jersey.

第二章 误差分析

数值代数中提供的数值解法，通常需要在高速电子计算机上进行。代数四则运算通常都带有舍入误差，原始数据也可能有误差，因此，求得的代数问题的解往往是近似解，近似解的精度不但与问题的性态有关，而且与算法也有关。本章将利用“向后误差分析”的方法来讨论方程组的性态问题和算法稳定性等有关概念。

§1 方程组的性态和算法的稳定性

1.1 方程组的性态

用数值方法解具体问题，我们最关心的是解的精度如何？通过实际计算，我们发现，解的精度与原始数据的精度，数学模型和算法等都有关系。下面从数值问题的一般概念出发，逐步介绍与数值问题的解的精度有关的内容。

例2.1.1 线性方程组

$$\begin{bmatrix} 1 & 10^4 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 10^4 \\ 2 \end{bmatrix}$$

的精确解如按三位有效数字进行舍入为： $x_1 = x_2 = 1.00$ 。如果对原始数据和中间结果均取三位有效数字利用部分主元素法求解，得到

$$x_1 = 0.00, x_2 = 1.00;$$

利用全主元法求解，得到

$$x_1 = 1.00, x_2 = 1.00.$$

例2.1.2 矩阵

$$A = \begin{bmatrix} 1 & 10^4 \\ 1 & 1 \end{bmatrix}$$

如果对元素取三位有效数字，利用特征多项式求其特征值，得到

$$\lambda_1 = 101, \lambda_2 = -99.$$

一般地，我们可以把数值问题描述成从原始数据 $d \in D$ 到数值结果 $s \in S$ 的映射 ϕ ，即

$$\phi: D \rightarrow S, \quad (2.1.1)$$

其中 D 表示原始数据 d 的集合， S 表示数值结果 S 的集合。这里， ϕ 包括数值问题的描述和解法等双重意义， ϕ 的自变量 d 是所论问题的原始输入数据，而 $\phi(d)$ (即 S) 就是问题的解。

在例 2.1.1 中

$$d = \begin{bmatrix} 1 & 10^4 & 10^4 \\ 1 & 1 & 2 \end{bmatrix}.$$

我们可分别用 $S = \phi(d)$ ， $S_1 = \phi_1(d)$ ， $S_2 = \phi_2(d)$ 表示相应问题的方程组的精确解，用部分主元素法求得的解和用全主元素法求得的解。此时

$$S = \phi(d) = \begin{bmatrix} 10000 \\ 9999 \\ 9999 \\ 9999 \end{bmatrix}, \quad S_1 = \phi_1(d) = \begin{bmatrix} 0.00 \\ 1.00 \end{bmatrix},$$

$$S_2 = \phi_2(d) = \begin{bmatrix} 1.00 \\ 1.00 \end{bmatrix}.$$

将某个问题化成数值问题时，我们往往希望，所求的解 $\phi(d)$ 存在而且唯一。如果用非负数 $\|\phi(d + \delta d) - \phi(d)\|$ 来表示用两组数据 d 和 $d + \delta d$ 所求的相应解的差异，当

(1) 存在 $\varepsilon > 0$ ，使对所有满足 $\|\delta d\|_D < \varepsilon$ 的 δd ，解 $\phi(d + \delta d)$ 均存在而且唯一，

(2) 当 $\|\delta d\|_D \rightarrow 0$ 时, $\|\phi(d + \delta d) - \phi(d)\|_S \rightarrow 0$, 则称该问题对给定的数据 d 是适定的. 其中 $\|\cdot\|_D, \|\cdot\|_S$ 分别表示定义在 D 和 S 上的某种范数. 条件(2)表示解对原始数据的连续依赖性.

如果问题的解多于一个, 或者解不连续依赖于原始数据, 那么, 应用数值方法求解时, 就会产生求得的解是否可信的问题, 因此, 本书只限于讨论适定的问题.

在实际问题中, 我们往往只知道原始数据的近似值 \bar{d} , 对于确定的 ϕ , 一般只能计算 $\phi(\bar{d})$. 如果 $\|d - \bar{d}\|_D \leq \delta$, 而 $\|\phi(\bar{d}) - \phi(d)\|_S \leq \varepsilon$, 其中 ε 可以估计成较小的正数, 我们就称该计算问题是“良态”的. 然而, 如 $\|\phi(\bar{d}) - \phi(d)\|$ 很大, 我们就称该计算问题是“病态”的. 对于病态问题, 当原始数据有微小的变化时, 可能引起数值结果有很大的波动. 这时若要求数值结果达到某种精度, 就必须要求原始数据具有更高的精度才行. 有时, 对于固定的 ϕ 和某种运算工具(如电子计算机), 要求数值结果达到指定的精度甚至是不可能的. 值得注意的是: 对于同一组原始数据 d , 关于不同的 ϕ , 其性态也可能不同.

例2.1.3 线性方程组

$$\begin{bmatrix} 3.000 & 4.127 \\ 1.000 & 1.374 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 15.41 \\ 5.147 \end{bmatrix} \quad (2.1.2)$$

有解 $x_1 = 13.6658, x_2 = -6.2$

$$\text{但线性方程组 } \begin{bmatrix} 3.000 & 4.122 \\ 1.000 & 1.374 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 15.41 \\ 5.147 \end{bmatrix} \quad (2.1.3)$$

没有解. 方程组(2.1.2)与(2.1.3)仅在系数矩阵的元素(1,2)上有0.005的差异, 但前者有解, 后者却没有解. 这是由于方程组(2.1.2)是“病态”所引起的. 方程组的病态也可以说成系数矩阵关于方程组求解问题是“病态”的.

例2.1.4 矩阵

$$A = \begin{bmatrix} 3.000 & 4.127 \\ 1.000 & 1.374 \end{bmatrix}$$

的特征值是 $\lambda_1 = -0.001$, $\lambda_2 = 4.375$, 而矩阵

$$B = \begin{bmatrix} 3.000 & 4.122 \\ 1.000 & 1.374 \end{bmatrix}$$

的特征值为 $\lambda_1 = 0$, $\lambda_2 = 4.374$.

在例2.1.4中, 矩阵 A 与 B 在元素 $(1, 2)$ 上同样也有 0.005 的差异, 但是它们所求得特征值仅有 0.001 的差别, 这是由于矩阵 A 的特征值问题是“良态”的. 这两个例子说明了问题的性态与问题的描述有关. 关于方程组的性态问题, 在国内外已有很多研究, 近年来, 随着电子计算机的计算速度和存贮量的提高, 不少原来计算中遇到的难题已得到了解决. 因此, 愈来愈多的人对研究“病态”问题发生了兴趣, 并且取得了不少研究成果.

本节着重研究线性方程组的性态, 特别是要讨论如何从数值上去判定线性方程组的性态等问题.

首先我们来看一下方程组 (2.1.2). 根据高等代数的知识可知, 线性方程组的解与系数行列式密切相关. 该方程组的系数行列式为

$$D = \begin{vmatrix} 3.000 & 4.127 \\ 1.000 & 1.374 \end{vmatrix} = -0.006,$$

其绝对值较小. 那么, 是否系数行列式的绝对值小, 方程组就是病态的呢? 下面先看两个例题.

例2.1.5 线性方程组

$$\begin{cases} 0.221x_1 + 0.054x_2 = 0.329 \\ 0.049x_1 + 0.034x_2 = 0.117 \end{cases}; \quad (2.1.4)$$

的系数行列式为

$$D = \begin{vmatrix} 0.221 & 0.054 \\ 0.049 & 0.034 \end{vmatrix} = 0.004868,$$

其解为

$$x_1 = 1, x_2 = 2.$$

线性方程组

$$\begin{cases} 0.221x_1 + 0.059x_2 = 0.329; \\ 0.049x_1 + 0.034x_2 = 0.117 \end{cases}$$

的解为

$$x_1 = 0.93, x_2 = 2.1$$

从例 2.1.5 可以看出, 系数行列式的绝对值小, 方程组并不一定病态。那么, 用什么数值来“度量”方程组的性态才合理呢? 方程组 (2.1.2) 与 (2.1.4), 虽然系数行列式的绝对值都较小, 但是两者却有很大不同, 前者的系数矩阵的列几乎线性相关, 而后者却没有这一特点。方程组 (2.1.2) 可以改写成

$$\begin{pmatrix} 3.000 \\ 1.000 \end{pmatrix} x_1 + \begin{pmatrix} 4.127 \\ 1.374 \end{pmatrix} x_2 = \begin{pmatrix} 15.41 \\ 5.147 \end{pmatrix}.$$

把 x_1 和 x_2 的系数及常数项记作向量 a_1, a_2, b , 则方程组 (2.1.2) 可以写成

$$a_1 x_1 + a_2 x_2 = b.$$

今把向量 a_1, a_2, b 标准化, 记 $\bar{a}_1 = \frac{a_1}{\|a_1\|}, \bar{a}_2 = \frac{a_2}{\|a_2\|}, \bar{b} = \frac{b}{\|b\|}$ (则 $\|\bar{a}_1\| = \|\bar{a}_2\| = \|\bar{b}\| = 1$), 组成新方程组

$$\bar{a}_1 y_1 + \bar{a}_2 y_2 = \bar{b}, \quad (2.1.5)$$

其中,

$$\frac{\|a_1\|}{\|b\|} x_1 = y_1, \quad \frac{\|a_2\|}{\|b\|} x_2 = y_2.$$

倘若(2.1.5)的解 y_1, y_2 已经求出, 并把关系式

$$\overline{b} - \overline{a}_2 y_2 = \overline{a}_1 y_1$$

画在平面坐标系 $O-xy$ 上, 如图2.1.

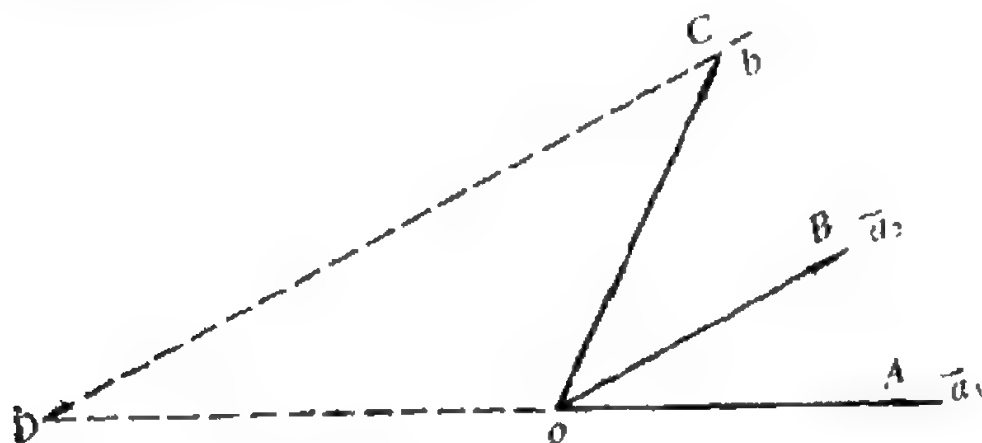


图2.1 列向量的变化对解的影响

其中 $CD // \overline{a}_2$, OD 是 AO 的延长线, 故有

$$\overline{b} - \overline{a}_2 y_2 = \overline{a}_1 y_1,$$

从而方程(2.1.5)的解的绝对值 $|y_1|$ 和 $|y_2|$ 分别是线段 DO 和 CD 的长. 从图2.1中可以清楚地看到, 当 \overline{a}_1 和 \overline{a}_2 几乎线性相关时, 即向量 \overline{a}_1 和 \overline{a}_2 的夹角很小时, \overline{a}_1 和 \overline{a}_2 稍有变动, 则向量 DO 和 CD 的长就可能发生很大的变化, 因而方程组是病态的.

上面我们从直观上说明了系数矩阵的列向量的线性关系与方程组的性态有着密切的关系. 把图形中的线性关系变成数值关系, 从而去判定方程组的性态, 将是很必要的.

设线性方程组为

$$Ax = b,$$

其中

$A = (\overline{a}_1, \overline{a}_2, \dots, \overline{a}_n)$, \overline{a}_i 是 A 的第 i 个列向量, $x = (x_1, x_2, \dots, x_n)^T$, $b = (\overline{a}_1, \overline{a}_2, \dots, \overline{a}_n)^T$. 根据线性代数知识, 列向量系 $\{\overline{a}_1, \overline{a}_2, \dots, \overline{a}_n\}$ 为线性相关的充要条件是, 该向量系的Gram行列式

$G[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n] = \det(\mathbf{a}_i^T \cdot \mathbf{a}_j) = 0$. Gram 行列式的几何意义是: 以向量系 $\{\mathbf{a}_1, \dots, \mathbf{a}_n\}$ 中各列向量为边所组成的平行多面体的体积的平方, 若向量系线性相关, 则该体积为零. 显然, 向量系中的非零向量是否线性相关与它们本身的长度是无关的, 但是体积的大小却与向量的长度有关. 为了使 $G[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$ 能表明向量系的相关程度, 以后, 我们可以将向量系中的向量标准化, 即使 $\|\mathbf{a}_i\| = 1 (i = 1, 2, \dots, n)$. 若

$$G[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n] \leq \varepsilon,$$

则称向量系 $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ 为 ε -线性相关.

据第一章§3的知识, 对于实非奇异矩阵 $A = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)$ 必存在正交矩阵 Q 使

$$A = QR,$$

其中

$$R = \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ 0 & r_{22} & \dots & r_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & r_{nn} \end{bmatrix},$$

$r_{ii} > 0 (i = 1, 2, \dots, n)$, 而且分解是唯一的. 因此

$$A^T A = (QR)^T (QR) = R^T R,$$

从而

$$\begin{aligned} G[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n] &= \det(A^T A) = (\det(R))^2 \\ &= r_{11}^2 \cdot r_{22}^2 \cdot \dots \cdot r_{nn}^2. \end{aligned}$$

由于列向量 \mathbf{a}_i 均已单位化, 所以 $|r_{ii}| \leq 1, i = 1, 2, \dots, n$. 设

$$p_1 = \text{Max}(|r_{11}|, |r_{22}|, \dots, |r_{nn}|);$$

$$p_n = \text{Min}(|r_{11}|, |r_{22}|, \dots, |r_{nn}|).$$

因此, p_n 愈小, 向量系愈接近相关.

定义2.1.1 称 $k^{(1)}(A) = \frac{1}{p_n}$ 为矩阵 A 的一阶条件数。

当 p_n 很小或 $k^{(1)}(A)$ 很大时，向量系就接近线性相关，并称方程组是病态的。

定义2.1.1中所述矩阵 A 的一阶条件数 $k^{(1)}(A)$ 与数值分析中定义的矩阵 A 的条件数 $\text{Cond}(A) = \|A\| \cdot \|A^{-1}\|$ 虽然形式不同，但它们却是从不同角度刻划了线性方程组的性态，在表达形式和使用方面也各具特色。它们之间还有如下关系：

定理2.1.1 设 $k^{(1)}(A)$ 是矩阵 A 的一阶条件数， $\text{Cond}(A) = \|A\| \|A^{-1}\|$ 是通常的条件数，则

$$K^{(1)}(A) \leq \text{Cond}(A) \leq \sqrt{n} \cdot 2^{n-1} K^{(1)}(A). \quad (2.1.6)$$

这个定理的证明，可以参看高等学校计算数学学报1979年第一期何旭初先生的文章。

为了说明 $K^{(1)}(A)$ 的大小能刻划方程组的性态，只需说明初始数据有微小变化时， $K^{(1)}(A)$ 的大小对方程组的解有多大影响就行了。根据定理2.1.1，只需对 $\text{Cond}(A)$ 进行这样的说明即可。为此，设方程组 $Ax = b$ 中， A 和 b 分别有误差 δA 和 δb ，从而解有误差 δx ，则它们之间有如下关系：

定理2.1.2 设在 $Ax = b \neq 0$ 中， A 是非奇异阵，且

$$(A + \delta A)(x + \delta x) = b + \delta b.$$

如果 $\|A^{-1}\| \|\delta A\| < 1$ ，则

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{Cond}(A)}{1 - r \text{Cond}(A)} \left(r + \frac{\|\delta b\|}{\|b\|} \right), \quad (2.1.7)$$

其中

$$r = \frac{\|\delta A\|}{\|A\|}.$$

证明 由 $(A + \delta A)(x + \delta x) = b + \delta b$

可得

$$(A + \delta A)\delta x = \delta b - (\delta A)x,$$

从而有

$$A(I + A^{-1}\delta A)\delta x = \delta b - (\delta A)x. \quad (2.1.8)$$

因为

$$\|A^{-1}\delta A\| \leq \|A^{-1}\|\|\delta A\| < 1,$$

所以 $(I + A^{-1}\delta A)^{-1}$ 存在, 于是由 (2.1.8) 得到

$$\delta x = (I + A^{-1}\delta A)^{-1}A^{-1}(\delta b - \delta A \cdot x),$$

故

$$\|\delta x\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\|\|\delta A\|} (\|\delta b\| + \|\delta A\|\|x\|),$$

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\|\|\delta A\|} \left(\frac{\|\delta b\|}{\|x\|} + \|\delta A\| \right).$$

又因为

$$\|A\|\|x\| \geq \|b\|,$$

所以

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\|\|A\|}{1 - \|A^{-1}\|\|\delta A\|} \left(\frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right),$$

即

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{Cond}(A)}{1 - r \cdot \text{Cond}(A)} \left(r + \frac{\|\delta b\|}{\|b\|} \right).$$

推论 2.1.1 在定理 2.1.2 中, 如果 $\delta b = 0$, 则

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{r \cdot \text{Cond}(A)}{1 - r \cdot \text{Cond}(A)}, \quad (2.1.9)$$

如果 $\delta A = 0$, 则

$$\frac{\|\delta x\|}{\|x\|} \leq \text{Cond}(A) \frac{\|\delta b\|}{\|b\|}. \quad (2.1.10)$$

从(2.1.7), (2.1.9), (2.1.10)可以看到, 如果 $\text{cond}(A)$ 很大, 那末, 解的相对误差界将比 $\|\delta A\|/\|A\|$ 和 $\|\delta b\|/\|b\|$ 大很多, 有可能使方程组的计算解与精确解相差很大, 即原始数据有较小的变化时, 方程组的解可能有很大的变化, 也就是方程组可能是病态的。

关于矩阵性态的定量的研究, 还在继续发展, 怎样用便于估计的量去全面刻画线性方程组的性态问题, 引起了不少人的兴趣。目前发表的论条件数的各种定义, 正是反映了这方面的努力。

1.2 算法的稳定性概念

数学问题在求解过程中, 往往需要进行成千上万次的运算, 由于舍入的原因, 每一步运算都会带来误差, 这种舍入误差的积累是否会影响解的精度呢? 这就要看算法的“好坏”了。在数学上, 衡量算法“好坏”的主要标志之一是在看算法是否稳定。从输入数据中所含的误差对解的精度影响以及从输出结果来推测输入数据中的误差来分析, 有两种稳定性的定义, 即向前稳定与向后稳定的定义值得提出来。

“向前稳定”是从“向前误差分析”所得到的结果来衡量的。所谓向前误差分析, 是指对每一步计算, 都去找出精确值与计算值之间的误差界 (此时, 输入数据的误差界被认为是设定的), 随着计算过程的向前推移而逐步向前分析, 直到得出最后计算结果与精确结果之间的误差界为止。如果这个界限比较小, 或者在所“允许”的范围之内, 我们就称这个算法是稳定的, 更确切地说, 就叫做“向前稳定”。

向前误差分析在使用时, 由于精确值事先不知道, 每步计算都要估计误差限, 不但工作量大, 而且容易产生过估现象。因此, 使用向前误差分析去分析算法的稳定性时, 历史上曾出

现过对某些算法的悲观论调。在很长一段时间内，研究误差分析的进展不大。直到提出“向后误差分析”之后，误差的分析研究才获得了比较全面的进展。

“向后稳定”是由向后误差分析的结果来判定的。设以 d 表示某数学问题的输入数据， $d \in D$ ，用 $\phi(d)$ 表示该问题的精确解，即求解的每一步中的中间结果都是精确地进行的，不带舍入误差的，而以 $\phi^*(d)$ 表示同一问题的计算解，即求解的每一步都是考虑了舍入误差的。如果把 $\phi^*(d)$ 写成 $\phi^*(d) = \phi(\bar{d})$ 的形式， $\bar{d} \in D$ ，并且 d 与 \bar{d} 的差别即 $\|d - \bar{d}\|$ 也不太大时，我们就称算法 ϕ^* 关于 D 是稳定的。可以把 \bar{d} 看成是 d 的带有误差的输入数据。由于原始数据一般总有误差，因此，这样处理是合于实际的。把计算过程的误差返回到原始数据的误差，这样的误差分析方法通常称为向后误差分析方法，如果算法经这样的分析是稳定的，则称之为是向后稳定的。

用一个稳定算法去解良态问题时，由于计算解 $\phi^*(d)$ 必须接近某个精确解 $\phi(\bar{d})$ ，其中 \bar{d} 接近于 d 。又由于问题是良态的，即当 \bar{d} 接近 d 时， $\phi(\bar{d})$ 接近 $\phi(d)$ ，因此 $\phi^*(d)$ 也接近于 $\phi(d)$ 。也就是说，用一个稳定算法去解良态问题时，计算解必定接近于精确解。如果用一个稳定算法去解病态问题，那末，由于 \bar{d} 接近于 d 时， $\phi(\bar{d})$ 不一定接近于 $\phi(d)$ ，从而 $\phi^*(d)$ 不一定接近于 $\phi(d)$ ，也就是说，稳定算法解病态问题时，得到的计算解与精确解可能差别很大。对于病态问题的求解，不是一般稳定算法所能解决的，必须要用特殊的算法来处理。

§2 误差分析

前一节我们讨论了数学问题的性态和算法的稳定性，这两个问题都是由于数据或运算中出现误差而引起的。为了深入地

讨论这些问题，必须进行误差分析。在四十年代以前，几乎没有人系统地研究过这个问题。直到四十年代中期，由于电子计算机的出现，才引起人们的巨大兴趣，开始的一些分析，然而，结论是很悲观的。例如，认为消去法不能用于解高阶方程组问题等。到1947年 Von Neumann 等的文章才澄清了这种悲观的论调，并含蓄地引进了“向后误差分析”的思想。直到1954年，Givens的文章才明确了“向后误差分析”的方法。在 J. H. Wilkinson 写的 “The Algebraic Eigenvalue Problem” 一书中，应用这种方法系统地分析了数值代数中众多的具体的算法，得到相当满意的结果。限于篇幅，本节将主要使用向后误差分析的方法，对代数中的基本运算和消去法进行浮点误差分析。

2.1 浮点基本运算的误差分析

对于浮点计算的各种模式，我们不能一一介绍，这里仅采用一种规格化的浮点计数系统来说明舍入误差分析方法。

对于任何实数 x ，通常都可找到一个整数 a 及一个小数 b ，使

$$x = 2^a \cdot b, \quad (2.2.1)$$

为了使这种表示形式是唯一确定的，可以规定

$$x \neq 0 \text{ 时, } \frac{1}{2} \leq |b| < 1, \quad (2.2.2)$$

$$x = 0 \text{ 时, } a = b = 0. \quad (2.2.3)$$

数 b 通常被表示成二进制形式，它可以是有限位的也可以是无限位的。在电子计算机上通常用若干位二进制数位，例如用 t 位二进制数位表示 b ，当 b 的二进制数位多于 t 位时，就要在 b 的第 $(t+1)$ 位上进行舍入，即当 b 的第 $(t+1)$ 位为 1 时，就在 b 的第 t 位上进 1，同时把 b 的从 $(t+1)$ 位以后的数甩掉，而

当 b 的第 $(t+1)$ 位是零时，就只甩掉 $(t+1)$ 位以后的数，而不再在 t 位上进位。这样得出的数的表示形式，叫做 x 的浮点规格化的机器表示形式，简称为机器表示形式。 a 叫做 x 的阶码， b 叫做尾数。由此可见数 x 与其机器表示的值可能并不精确地相等，如果有误差，则造成误差的原因是“舍入”，而误差的大小取决于 t —— 机器的字长。即

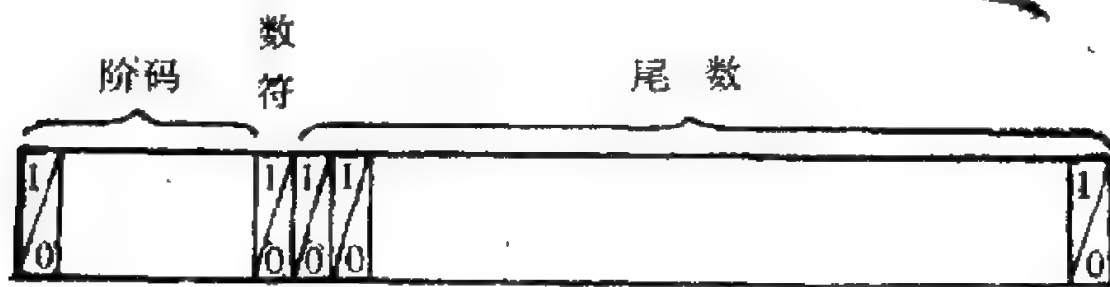


图2.2 数的机器表示

如果用 $fl(x)$ 来表示 x 的机器表示值，则含于 $fl(x)$ 中的绝对误差 ε_e 和相对误差 ε_r 应满足

$$fl(x) = x + \varepsilon_e, \quad |\varepsilon_e| \leq 2^{a-t-1},$$

$$fl(x) = x + x \cdot \frac{\varepsilon_e}{x} = x \left(1 + \frac{\varepsilon_e}{x} \right) = x(1 + \varepsilon_r),$$

$$|\varepsilon_r| = \left| \frac{\varepsilon_e}{x} \right| \leq \frac{2^{a-t-1}}{|2^a \cdot b|} \leq 2^{-t}.$$

下面我们讨论算术运算的舍入误差，由于不可能同时研究所有的计算机，在此假定计算机的累加器是双字长的，即可以接受双字长的数。

设 x, y 为规格化的浮点数，并且用

$$fl(x+y), fl(x-y), fl(x \times y), fl(x/y)$$

分别表示进行浮点加，减，乘，除四则运算后的结果。

2.1.1 加法(减法)

设 $x = 2^{b_1}a_1, y = 2^{b_2}a_2, b_1 > b_2$ ；求 $fl(x \pm y)$ 。

如果 $b_1 - b_2 \leq t$, 则将 a_2 右移 $b_1 - b_2$ 位 (将移出去的 $b_1 - b_2$ 位保留在累加器的后半部), 使 x_1 与 x_2 同阶, 然后求 $a_1 \pm a_2 \times 2^{b_1 - b_2}$, 并进行规格化, 对前 t 位进行舍入 (即在 $t+1$ 位上加上“ ± 1 ”), 最后再调整阶码。其中只有一步舍入。设最后结果为 $2^{b_3} a_3$, 则

$$fl(x \pm y) = x \pm y + e_r, \quad |e_r| \leq 2^{b_3 - t - 1},$$

$$\begin{aligned} fl(x \pm y) &= (x \pm y) + (x \pm y) \frac{e_r}{x \pm y} \\ &= (x \pm y) \left(1 + \frac{e_r}{x \pm y} \right) = (x \pm y) (1 + \varepsilon_r), \end{aligned}$$

$$|\varepsilon_r| = \frac{|e_r|}{|x \pm y|} \leq \frac{2^{b_3 - t - 1}}{|2^{b_3} a_3|} \leq \frac{2^{b_3 - t - 1}}{2^{b_3} - 1} = 2^{-t}.$$

于是可以写成如下形式

$$\begin{cases} fl(x \pm y) = (x \pm y) (1 + \varepsilon_r), \\ |\varepsilon_r| \leq 2^{-t} \quad (b_1 \geq b_2). \end{cases} \quad (2.2.4)$$

当 $b_1 \leq b_2$ 时, (2.2.4) 同样满足。(2.2.4) 表明 x 与 y 的计算和 (或差) 是 $x(1 + \varepsilon_r)$ 与 $y(1 + \varepsilon_r)$ 的精确和 (或差), 这里 $|\varepsilon_r| \leq 2^{-t}$ 。

2.1.2 乘法

设 $x = 2^{b_1} \cdot a_1$, $y = 2^{b_2} \cdot a_2$ 求 $fl(x \times y)$ 。

首先是指数相加: $b_1 + b_2$, 然后精确求出 $a_1 \cdot a_2$ 的双字长, 即 $2t$ 位的结果, 再规格化, 调整阶码, 最后将 $2t$ 位字长的数舍入成 t 位。若 x 或 y 为零, 则乘积直接取零。由此可见也只有最后一步舍入。设其结果为 $2^{b_4} \cdot a_4$, 则

$$fl(x \times y) = x \times y + e_r, \quad |e_r| \leq 2^{b_4 - t - 1},$$

$$fl(x \times y) = (x \times y) \left(1 + \frac{e_r}{x \times y} \right) = (x \times y) (1 + \varepsilon_r),$$

$$|\varepsilon_r| = \frac{|e_r|}{|x \times y|} \leq \frac{2^{b_4 - t - 1}}{|2^{b_4} \cdot a_4|} \leq 2^{-t},$$

从而有

$$\begin{cases} fl(x \times y) \equiv (x \times y)(1 + \varepsilon_r), \\ |\varepsilon_r| \leq 2^{-t}, \end{cases} \quad (2.2.5)$$

即 x 与 y 的计算积是 $x(1 + \varepsilon_r)$ 与 y 或 x 与 $y(1 + \varepsilon_r)$ 或 $x(1 + \varepsilon_r)^{\frac{1}{2}}$ 与 $y(1 + \varepsilon_r)^{\frac{1}{2}}$ 的精确积, 其中 $|\varepsilon_r| \leq 2^{-t}$. 这三种形式是可以视分析的方便而任选的.

例2.2.1 设 $x = 2^{101} \cdot 0.101101$ $y = 2^{11} \cdot 0.111101$ 求 $fl(x - y)$ 和 $fl(x \times y)$.

解 (1) 因为 $t = 110$ 且 $b_1 - b_2 = 101 - 11 = 10 \leq 110$

所以 $\bar{a}_3 = a_1 - a_2 \times 2^{-10} = 0.011110$,

将 $\bar{a}_3 = 0.011110$ 规格化得到 0.111100 .

对 0.111100 进行舍入后得 $a_3 = 0.111100$.

调整阶码算得最后结果为

$$fl(x - y) = 2^{100} \cdot 0.111100 = (x - y)(1 + \varepsilon_r),$$

$$\text{其中 } |\varepsilon_r| \leq 2^{-110}.$$

(2) 因为 $a_1 \cdot a_2 = 0.101011001001$,

将 $a_1 \cdot a_2$ 舍入成 6 位, 得到 $a_4 = 0.101011$,

$$2^{b_1 + t_1} = 2^{101 + 11} = 2^{1000} = 2^{b_4},$$

所以

$$\begin{cases} fl(x \times y) = 2^{1000} \cdot 0.101011 = (x \times y)(1 + \varepsilon_r), \\ |\varepsilon_r| \leq 2^{-110}. \end{cases}$$

2.1.3 除法

求 $fl(x/y)$, 只要 $y \neq 0$, 同加法和乘法一样, 可以得到

$$\begin{aligned} fl(x/y) &\equiv (x/y)(1 + \varepsilon_r) \\ |\varepsilon_r| &\leq 2^{-t} \end{aligned} \quad (2.2.6)$$

即 x 与 y 的计算商是 $x(1 + \varepsilon_r)$ 与 y 或 x 与 $y/(1 + \varepsilon_r)$ 的精确商,

其中 $|e_r| \leq 2^{-t}$.

综合上面的讨论, 如果用 θ 表示加, 减, 乘, 除四则运算中的任何一种, 则有下列统一的表示:

$$\begin{cases} fl(x\theta y) = (x\theta y)(1 + e_r), \\ |e_r| \leq 2^{-t}. \end{cases} \quad (2.2.7)$$

若设

$$1 + e_r = \frac{1}{1 + e'_r}, \text{ 其中 } |e_r| \leq 2^{-t},$$

则

$$1 + e'_r = \frac{1}{1 + e_r} \approx 1 - e_r + e_r^2 - \dots + \dots,$$

即

$$e'_r = -e_r + e_r^2 - e_r^3 + \dots - \dots.$$

当 e_r 很小时, 可以忽略 e_r 的高次项, 于是

$$e'_r \approx -e_r,$$

从而(2.2.7)还可以写成如下适用的形式

$$\begin{cases} fl(x\theta y) = (x\theta y) / (1 + e'_r), \\ |e'_r| \leq 2^{-t}. \end{cases} \quad (2.2.8)$$

2.1.4 连加和连乘

设 x_1, x_2, \dots, x_n 均为标准浮点数, 求 $fl(x_1 + x_2 + \dots + x_n)$ 和 $fl(x_1 x_2 \dots x_n)$. 计算从左向右进行, 每进行一次都进行舍入, 然后再做下一次.

(1) 连加

设

$$\begin{cases} s_1 = x_1, \\ s_2 = fl(s_1 + x_2) \equiv (s_1 + x_2)(1 + e_2), \\ s_k = fl(s_{k-1} + x_k) \equiv (s_{k-1} + x_k)(1 + e_k), \\ \quad k = 3, 4, \dots, n \\ |e_k| \leq 2^{-t} \quad (k = 2, 3, \dots, n), \end{cases} \quad (2.2.9)$$

则有

$$\begin{cases} s_n = fl(x_1 + x_2 + \dots + x_n) \equiv x_1(1 + \eta_1) + x_2(1 + \eta_2) \\ \quad + \dots + x_n(1 + \eta_n), \\ (1 - 2^{-t})^{n+1-k} \leq 1 + \eta_k \leq (1 + 2^{-t})^{n+1-k}, \\ k = 1, 2, \dots, n. \end{cases} \quad (2.2.10)$$

证明 将(2.2.9)中的 s_k 逐项代入, 可得

$$s_n = x_1(1 + \eta_1) + x_2(1 + \eta_2) + \dots + x_n(1 + \eta_n).$$

其中

$$1 + \eta_1 = 1 + \eta_2 = (1 + \varepsilon_2)(1 + \varepsilon_1) \dots (1 + \varepsilon_n),$$

$$1 + \eta_k = (1 + \varepsilon_k)(1 + \varepsilon_{k+1}) \dots (1 + \varepsilon_n),$$

$$k = 3, 4, \dots, n.$$

从而有

$$(1 - 2^{-t})^{n-k+1} \leq 1 + \eta_k \leq (1 + 2^{-t})^{n-k+1}.$$

由于 $(1 \pm 2^{-t})^k$ 的形式经常出现, 所以引入简化条件. 设

$$k \cdot 2^{-t} < 0.1 \quad (2.2.11)$$

此条件并不苛刻, 例如当 $t = 40$ 时, k 可达 10^{11} .

由于

$$\begin{aligned} (1 + 2^{-t})^k &\leq 1 + k \cdot 2^{-t} + \frac{k(k-1)}{2!} (2^{-t})^2 + \dots \\ &= 1 + k \cdot 2^{-t} \left(1 + \frac{k-1}{2!} (2^{-t}) + \frac{(k-1)(k-2)}{3!} (2^{-t})^2 + \dots \right) \\ &< 1 + k \cdot 2^{-t} \left(1 + \frac{1}{2!} (0.1) + \frac{1}{3!} (0.1)^2 + \dots \right) \\ &= 1 + k \cdot 2^{-t} \cdot \frac{1}{0.1} (e^{0.1} - 1) < 1 + (1.06) k \cdot 2^{-t}. \end{aligned}$$

类似可得

$$(1 - 2^{-t})^k \geq 1 - (1.06) k \cdot 2^{-t}.$$

令

$$(1.06) 2^{-t} = 2^{-t_1},$$

则

$$t_1 = t - \log_2(1.06) \doteq t - 0.084.$$

于是 $(1 + 2^{-t})^k < 1 + k \cdot 2^{-t}, (1 - 2^{-t})^k > 1 - k \cdot 2^{-t}$

其中 $t_1 = t - 0.084$

所以 (2.2.10) 的第 2 式可写成

$$\begin{aligned} 1 - (1.06)(n - k + 1)2^{-t} &< 1 + \eta_k < 1 + \\ &(1.06)(n - k + 1)2^{-t}, \quad (2.2.12) \\ |\eta_k| &< (1.06)(n - k + 1)2^{-t} = (n - k + 1)2^{-t_1}. \end{aligned}$$

综上所述可得

$$\begin{cases} s_n = f(x_1 + x_2 + \cdots + x_n) = x_1(1 + \eta_1) + \cdots \\ \quad + x_n(1 + \eta_n), \\ |\eta_k| \leq (n - k + 1) \cdot 2^{-t_1} \quad (k = 1, 2, \cdots, n), \end{cases} \quad (2.2.13)$$

(2.2.10), (2.2.13) 表明, 连加的计算和是带有相对误差的各数 $x_k(1 + \eta_k)$ 的精确和, 诸相对误差限的大小, 因参加运算的各数的先后次序而异, 其先参加运算者在和式中引起的误差的增长也较大. 因此, 求和时应先安排数值较小的数参加运算, 这样做, 在多数情况下能取得较高的精度.

(2) 连乘

设 $p_1 = x_1,$

$$p_2 = fl(p_1 x_2) \equiv (p_1 x_2)(1 + \varepsilon_2),$$

\vdots

$$p_k = fl(p_{k-1} x_k) \equiv (p_{k-1} x_k)(1 + \varepsilon_k),$$

$$|\varepsilon_k| \leq 2^{-t} \quad (k = 2, 3, \cdots, n),$$

则有

$$\begin{cases} p_n = fl(x_1 x_2 \cdots x_n) = x_1 x_2 \cdots x_n (1 + e), \\ (1 - 2^{-t})^{n-1} \leq 1 + e \leq (1 + 2^{-t})^{n-1}, \end{cases} \quad (2.2.14)$$

由于 $p_n = x_1 x_2 \cdots x_n (1 + \varepsilon_2)(1 + \varepsilon_3) \cdots (1 + \varepsilon_n)$, 且 $|\varepsilon_k| \leq 2^{-t}$,

若令

$$(1 + \varepsilon_2)(1 + \varepsilon_3) \cdots (1 + \varepsilon_n) = 1 + e,$$

则有

$$(1 - 2^{-t})^{n-1} \leq 1 + e \leq (1 + 2^{-t})^{n-1}. \quad (2.2.15)$$

其上、下界与连加不一样，它与计算次序无关。

2.1.5 内积

设 $x_i, y_i (i = 1, 2, \dots, n)$ 均为浮点规格化的数，求 $s_n = fl(x_1 y_1 + x_2 y_2 + \dots + x_n y_n)$ 。计算次序是先求出两两乘积的单字长浮点数，再依次相加，直到求出最后结果。即设

$$\begin{aligned} p_1 &= fl(x_1 y_1) = s_1, \\ p_2 &= fl(x_2 y_2), \quad s_2 = fl(s_1 + p_2), \\ &\vdots \\ p_k &= fl(x_k y_k), \quad s_k = fl(s_{k-1} + p_k), \\ &k = 2, 3, \dots, n, \end{aligned}$$

则有

$$\begin{cases} s_n = x_1 y_1 (1 + \varepsilon_1) + x_2 y_2 (1 + \varepsilon_2) + \dots + x_n y_n (1 + \varepsilon_n), \\ (1 - 2^{-t})^{n-k+2} \leq 1 + \varepsilon_k \leq (1 + 2^{-t})^{n-k+2}, \\ k = 1, 2, \dots, n, \end{cases} \quad (2.2.16)$$

因为 $p_k \equiv x_k y_k (1 + \xi_k), \quad |\xi_k| \leq 2^{-t},$

$$s_k \equiv (s_{k-1} + p_k) (1 + \eta_k), \quad |\eta_k| \leq 2^{-t},$$

所以

$$s_n = x_1 y_1 (1 + \varepsilon_1) + x_2 y_2 (1 + \varepsilon_2) + \dots + x_n y_n (1 + \varepsilon_n).$$

其中

$$\begin{aligned} 1 + \varepsilon_1 &= (1 + \xi_1) (1 + \eta_2) \dots (1 + \eta_n), \\ 1 + \varepsilon_2 &= (1 + \xi_2) (1 + \eta_3) \dots (1 + \eta_n), \\ &\vdots \\ 1 + \varepsilon_k &= (1 + \xi_k) (1 + \eta_k) \dots (1 + \eta_n), \\ &k = 3, 4, \dots, n. \end{aligned}$$

从而

$$\begin{cases} (1-2^{-t})^k \leq 1 + \varepsilon_1 \leq (1+2^{-t})^n, \\ (1-2^{-t})^{n-k+2} \leq 1 + \varepsilon_k \leq (1+2^{-t})^{n-k+2}, \\ k = 2, 3, \dots, n. \end{cases} \quad (2.2.17)$$

(2.2.17) 还可统一地写成

$$\begin{aligned} (1-2^{-t})^{n-k+2} \leq 1 + \varepsilon_k \leq (1+2^{-t})^{n-k+2} \\ k = 1, 2, \dots, n. \end{aligned} \quad (2.2.18)$$

当 $k \cdot 2^{-t} < 0.1$ 时, 据 (2.2.12) 有

$$1 - (1.06)(n-k+2)2^{-t} < 1 + \varepsilon_k < 1 + (1.06)(n-k+2)2^{-t},$$

即

$$|\varepsilon_k| < (n-k+2)(1.06)2^{-t} = (n-k+2)2^{-t_1}. \quad (2.2.19)$$

其中 $t_1 = t - 0.084$ 由此得到

$$\begin{cases} s_n = x_1 y_1 (1 + \varepsilon_1) + x_2 y_2 (1 + \varepsilon_2) + \dots \\ \quad + x_n y_n (1 + \varepsilon_n), \\ |\varepsilon_k| \leq (n-k+2)2^{-t_1} \quad k = 1, 2, \dots, n. \end{cases} \quad (2.2.20)$$

以上这些结果为算法的浮点误差分析及稳定性讨论提供了良好的基础。

2.2 消去法的浮点误差分析

利用前面的浮点误差分析的基本知识, 可以对求解线性方程组的消元法进行误差分析, 在此仅对列主元素法进行分析, 全主元法可以类似地进行。

在做误差分析时, 为叙述简单起见, 设消元是按自然次序进行的, 即每次消元时不需要进行行的交换, 按自然次序选取的主元就是列主元。设求解的线性方程组为

$$Ax = b.$$

首先考虑对 A 进行 LR 分解。设 $A = A^{(1)}$, 各步消元后产生的计算输出矩阵依次为 $A^{(2)}, A^{(3)}, \dots, A^{(n)} = \tilde{R} = (r_{ij})_{n \times n}$, 其中 $A^{(k)}$ 的元素为 $a_{ij}^{(k)}$, 而且当 $i < k$, $i > j$ 时 $a_{ij}^{(k)} = 0$, 即 $A^{(k)}$ 在

前 $(k-1)$ 列的对角元以下的元素均为0。由普通消去法对 $A^{(k)}$ 再进行一步消元后得 $A^{(k+1)}$ ，其 $a_{ij}^{(k+1)}$ 应为

$$a_{ij}^{(k+1)} = \begin{cases} 0, & i \geq k+1, j = k, \\ fl(a_{ij}^{(k)} - m_{ik} \cdot a_{kj}^{(k)}), & i \geq k+1, j \geq k+1, \\ a_{ij}^{(k)}, & \text{其他}, \end{cases} \quad (2.2.21)$$

其中 $m_{ik} = fl(a_{ik}^{(k)} / a_{kk}^{(k)}) \quad i \geq k+1$ 。

由于

$$m_{ik} = fl(a_{ik}^{(k)} / a_{kk}^{(k)}) = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} (1 + \varepsilon_{ik}), \quad |\varepsilon_{ik}| \leq 2^{-t}. \quad (2.2.22)$$

于是当 $i \geq k+1, j = k$ 时在 (i, k) 处的元素的精确结果应为

$$\begin{aligned} \bar{a}_{ik}^{(k+1)} &= a_{ik}^{(k)} - m_{ik} a_{kk}^{(k)} = a_{ik}^{(k)} - \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} (1 + \varepsilon_{ik}) a_{kk}^{(k)} \\ &= -a_{ik}^{(k)} \varepsilon_{ik}. \end{aligned} \quad (2.2.23)$$

这个结果一般不为零，当按消去法要求，强使该处的输出为零时，需要引入一个平衡的误差 $e_{ik}^{(k)}$ 使 $a_{ik}^{(k+1)} = \bar{a}_{ik}^{(k+1)} + e_{ik}^{(k)} = 0$ ，此时应有

$$e_{ik}^{(k)} = a_{ik}^{(k)} \varepsilon_{ik}. \quad (2.2.24)$$

当 $i \geq k+1, j \geq k+1$ 时

$$\begin{aligned} a_{ij}^{(k+1)} &= fl(a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)}) \\ &= [a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} (1 + \varepsilon_{ik})] (1 + \eta_{ij}), \\ |\varepsilon_{ij}| &\leq 2^{-t}, \quad |\eta_{ij}| \leq 2^{-t}. \end{aligned} \quad (2.2.25)$$

从而有

$$\begin{aligned} a_{ij}^{(k+1)} &= a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} + [a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} (1 + \varepsilon_{ik})] \eta_{ij} \\ &\quad - m_{ik} a_{kj}^{(k)} \varepsilon_{ij} \\ &= a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} + \frac{a_{ij}^{(k+1)}}{1 + \eta_{ij}} \eta_{ij} - m_{ik} a_{kj}^{(k)} \varepsilon_{ij}. \end{aligned} \quad (2.2.26)$$

改写(2.2.26)为

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} + e_{ij}^{(k)} \quad (2.2.27)$$

其中

$$e_{ij}^{(k)} = \frac{a_{ij}^{(k+1)}}{1 + \eta_{ij}} \eta_{ij} - m_{ik} a_{kj}^{(k)} e_{ij} \quad (i, j \geq k+1),$$

因为(2.2.21)的第三式不进行运算, 所以没有引进误差, 即

$$a_{ij}^{(k-1)} = a_{ij}^{(k)} + e_{ij}^{(k)}, \quad e_{ij}^{(k)} = 0. \quad (2.2.28)$$

综合(2.2.24)、(2.2.27)、(2.2.28)得到

$$a_{ij}^{(k+1)} = \begin{cases} a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} + e_{ij}^{(k)}, & i \geq k+1, j = k, \\ a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} + e_{ij}^{(k)}, & i \geq k+1, j \geq k+1, \\ a_{ij}^{(k)} + e_{ij}^{(k)}, & \text{其他,} \end{cases} \quad (2.2.29)$$

其中

$$e_{ij}^{(k)} = \begin{cases} a_{ik}^{(k)} e_{ik}, & |\varepsilon_{ik}| \leq 2^{-l}, i \geq k+1, j = k, \\ \frac{a_{ij}^{(k+1)}}{1 + \eta_{ij}} \eta_{ij} - m_{ik} a_{kj}^{(k)} e_{ij}, & |\eta_{ij}|, |\varepsilon_{ij}| \leq 2^{-l}, \\ i \geq k+1, j \geq k+1, \\ 0, & \text{其他.} \end{cases} \quad (2.2.30)$$

设 $E^{(k)}$ 是以 $e_{ij}^{(k)}$ 为元素的误差矩阵, 则(2.2.29)可以写成

$$A^{(k+1)} = A^{(k)} - M_k A^{(k)} + E^{(k)}, \quad (2.2.31)$$

其中

$$M_k = \begin{pmatrix} 0 & & & & \\ & \ddots & & & \\ & & 0 & & \\ & & & m_{k+1,k} & 0 \\ & & & & m_{k+2,k} & 0 \\ & & & & \vdots & \\ 0 & & & & m_{n,k} & \ddots & 0 \end{pmatrix}. \quad (2.2.32)$$

令 $k = 1, 2, \dots, n-1$, 相应地得到

$$A^{(2)} = A^{(1)} - M_1 A^{(1)} + E^{(1)},$$

$$A^{(3)} = A^{(2)} - M_2 A^{(2)} + E^{(2)},$$

\vdots

$$A^{(n)} = A^{(n-1)} - M_{n-1} A^{(n-1)} + E^{(n-1)}.$$

两边求和, 得到

$$A^{(n)} = A^{(1)} - \sum_{i=1}^{n-1} M_i A^{(i)} + \sum_{i=1}^{n-1} E^{(i)}. \quad (2.2.33)$$

注意 $A^{(k)}$ 的第 k 行与 $A^{(k+1)}, \dots, A^{(n)}$ 的第 k 行元素是一样的。而从 (2.2.32) 又可以看出, $M_k A^{(k)}$ 只与 $A^{(k)}$ 的第 k 行元素有关, 因此

$$M_k A^{(k)} = M_k A^{(n)} = M_k \tilde{R},$$

于是 (2.2.33) 可写成

$$\tilde{R} = A^{(1)} - \left(\sum_{i=1}^{n-1} M_i \right) \tilde{R} + \sum_{i=1}^{n-1} E^{(i)},$$

经整理得

$$(M_1 + M_2 + \dots + M_{n-1} + I) \tilde{R} = A + \sum_{i=1}^{n-1} E^{(i)}. \quad (2.2.34)$$

令

$$M_1 + M_2 + \dots + M_{n-1} + I = \tilde{L} = (l_{ij})_{n \times n},$$

$$E^{(1)} + E^{(2)} + \dots + E^{(n-1)} = E,$$

则

$$\tilde{L} \tilde{R} = A + E, \quad (2.2.35)$$

其中 \tilde{L} 是单位下三角阵, \tilde{R} 是上三角阵。从而我们得到了 A 的带有误差矩阵 E 的矩阵 $A + E$ 的精确的 LR 分解。

下面讨论 E 的元素的大小, 为此, 首先需要估计 E 的元素中的量 m_{ik} 和 $a_{ij}^{(k)}$ 的界。

由于假定按自然次序的消去法就是列主元消去法, 所以 $|a_{kk}^{(k)}| \geq |a_{ik}^{(k)}|$ ($i \geq k+1$), 且 $m_{ik} = fl(a_{ik}^{(k)}/a_{kk}^{(k)})$, 因此一般有

$$|m_{ik}| \leq 1, \text{ 即 } |l_{ik}| \leq 1.$$

令

$$\rho = \max_{i,j,k} |a_{ij}^{(k)}| / \|A\|_{\infty}, \quad (2.2.36)$$

则对任何 $1 \leq i, j, k \leq n$, 均有

$$|a_{ij}^{(k)}| \leq \rho \|A\|_{\infty} \quad (2.2.37)$$

及 $|r_{ij}| \leq \rho \|A\|_\infty$,

其中 r_{ij} 是上三角阵 $\tilde{\mathbf{R}}$ 的 (i, j) 元素。这里虽然不能对 ρ 的值做先验估计，但它是可以在消元过程中求出的。有了 m_{ik} 和 $a_{ij}^{(k)}$ 的界以后，就可以估计 $\varepsilon_{ij}^{(k)}$ 的界了。由 (2.2.30) 的第 1 式有

$$|\varepsilon_{ij}^{(k)}| = |a_{ik}^{(k)} \varepsilon_{ik}| \leq \rho \|A\|_\infty \cdot 2^{-i}, \\ i \geq k+1, j=k.$$

由 (2.2.30) 的第 2 式，得到

$$\begin{aligned} |\varepsilon_{ij}^{(k)}| &= \left| \frac{a_{ij}^{(k+1)} \eta_{ij}}{1 + \eta_{ij}} - m_{ik} a_{ij}^{(k)} \varepsilon_{ij} \right| \\ &\leq \frac{|a_{ij}^{(k+1)}| |\eta_{ij}|}{|1 + \eta_{ij}|} + |m_{ik}| |a_{ij}^{(k)}| |\varepsilon_{ij}| \\ &\leq \frac{\rho \|A\|_\infty \cdot 2^{-i}}{1 - 2^{-i}} + \rho \|A\|_\infty \cdot 2^{-i} \\ &\leq \rho \|A\|_\infty \cdot 2^{-i} \left(\frac{1}{1 - 2^{-i}} + 1 \right). \end{aligned}$$

由于

$$\frac{1}{1 - 2^{-i}} \leq 1 + 2^{-i},$$

所以

$$\begin{aligned} &1, i \geq k+1, j=k, \\ |\varepsilon_{ij}^{(k)}| &\leq \rho \|A\|_\infty \cdot 2^{-i} (2 + 2^{-i}), i \geq k+1, j \geq k+1, \quad (2.2.38) \\ &0, \text{其他.} \end{aligned}$$

于是得

$$|E^{(k)}| \leq \rho \|A\|_\infty \cdot 2^{-i} \begin{pmatrix} 0 & \vdots & 0 & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & \vdots & 1 & 2 + 2^{-i} & 2 + 2^{-i} \\ & \vdots & \vdots & \vdots & \vdots \\ & \vdots & 1 & 2 + 2^{-i} & 2 + 2^{-i} \end{pmatrix} \quad (\text{第 } k+1 \text{ 行})$$

又由 $E = \sum_{i=0}^{n-1} E^{(i)}$ 所以

$$|E| \leq \rho \|A\|_{\infty} \cdot 2^{-t} \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & 0 \\ 1 & 2 & 2 & \cdots & 2 & 2 \\ 1 & 3 & 4 & \cdots & 4 & 4 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & 3 & 5 & \cdots & 2n-4 & 2n-4 \\ 1 & 3 & 5 & \cdots & 2n-3 & 2n-2 \end{pmatrix} \\ + 2^{-t} \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 1 & 1 & \cdots & 1 & 1 \\ 0 & 1 & 2 & \cdots & 3 & 3 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 1 & 2 & \cdots & n-2 & n-2 \\ 0 & 1 & 2 & \cdots & n-2 & n-1 \end{pmatrix}.$$

$$\begin{aligned} \|E\|_{\infty} = \| |E| \|_{\infty} &\leq \rho \|A\|_{\infty} \cdot 2^{-t} \left\{ \left[\sum_{i=1}^n (2i-1-1) \right] + 2^{-t} \sum_{i=1}^{n-1} i \right\} \\ &\leq \rho \|A\|_{\infty} \cdot 2^{-t} (n^2 + 2^{-t} \cdot n^2/2) \\ &\leq \rho \|A\|_{\infty} \cdot n^2 \cdot 2^{-t+1}. \end{aligned}$$

定理2.2.1 使用列主元消去法时, 实际计算出来的三角因子 \tilde{L} 和 \tilde{R} 满足

$$\tilde{L} \tilde{R} = A + E, \quad (2.2.39)$$

$$\|E\|_{\infty} \leq n^2 \rho \|A\|_{\infty} 2^{-t+1}. \quad (2.2.40)$$

这个重要结果表明, 对 A 的实际计算所得到的三角分解, 可以看作是 $A + E$ 的精确的三角分解. 对于一般矩阵, 如果 ρ 不是很大, 那末, 在分解过程中由于舍入误差而引起的误差矩阵, 其元素的界与机器精度几乎是同阶的.

完成三角分解的误差分析以后, 对于求解方程组 $Ax = b$ 的

$$\tilde{L}y = b \quad \text{和} \quad \tilde{R}x = y, \quad (2.2.41)$$

因此,要完成对消元法的误差分析,首先要估计求解三角方程组

$$Ux = b \quad (2.2.42)$$

中的解的误差。不失一般性，设 U 为下三角阵，于是 (2.2.42) 可以写成

[illegible]

$$x_1 = fl(b_1/u_{11}) = b_1/(u_{11}(1 + \varepsilon_n)), \quad |\varepsilon_n| < 2^{-t}, \quad (2.2.44)$$

$$\begin{aligned} x_i &= fl\left(\frac{-u_{i1}x_1 - u_{i2}x_2 - \dots - u_{i,i-1}x_{i-1} + b_i}{u_{ii}}\right) \\ &= \frac{fl(-u_{i1}x_1 - u_{i2}x_2 - \dots - u_{i,i-1}x_{i-1} + b_i)}{u_{ii}(1 + \varepsilon_{ii})} \\ &= \frac{fl(-u_{i1}x_1 - u_{i2}x_2 - \dots - u_{i,i-1}x_{i-1}) + b_i}{u_{ii}(1 + \varepsilon_{ii})(1 + \varepsilon'_{ii})}. \end{aligned}$$

由 (2.2.20), (2.2.19) 得到

$$x_i = \frac{-u_{i1}(1 + \varepsilon_{i1})x_1 - u_{i2}(1 + \varepsilon_{i2})x_2 - \dots - u_{ii-1}(1 + \varepsilon_{ii-1})x_{i-1} + b_i}{u_{ii}(1 + \varepsilon_{ii})(1 + \varepsilon'_{ii})}, \quad (2.2.45)$$

其中

$$\left\{ \begin{aligned} |\varepsilon_{ii}|, |\varepsilon'_{ii}| &\leq 2^{-i}, \quad i = 1, 2, \dots, n, \\ |\varepsilon_{i1}| &< (i-1)(1.06)2^{-i}, \quad i = 2, 3, \dots, n, \\ |\varepsilon_{ik}| &< ((i-1) - k + 2)(1.06)2^{-i} \\ &= (i+1-k)(1.06)2^{-i} \\ &(i = 2, 3, \dots, n, \quad k = 2, 3, \dots, i-1). \end{aligned} \right. \quad (2.2.46)$$

由此可知

$$\begin{aligned} u_{11}(1 + \varepsilon_{11})x_1 &= b_1, \\ u_{i1}(1 + \varepsilon_{i1})x_1 + u_{i2}(1 + \varepsilon_{i2})x_2 + \cdots \\ &+ u_{i,i-1}(1 + \varepsilon_{i,i-1})x_{i-1} + u_{ii}(1 + \varepsilon_{ii})(1 + \varepsilon'_{ii})x_i = b_i, \\ (i &= 2, 3, \cdots, n). \end{aligned} \quad (2.2.47)$$

写成矩阵形式 $(U + \delta U)x = b$,

其中

$$\begin{aligned} \delta U &= \begin{pmatrix} u_{11}\varepsilon_{11} & & & \\ u_{21}\varepsilon_{21} & u_{22}\delta_{22} & & \\ \vdots & & & \\ u_{n1}\varepsilon_{n1} & u_{n2}\varepsilon_{n2} & \cdots & u_{nn-1}\varepsilon_{nn-1} & u_{nn}\delta_{nn} \end{pmatrix} \\ \delta_{ii} &= \varepsilon_{ii} + \varepsilon'_{ii} + \varepsilon_{ii}\varepsilon'_{ii}, \quad i = 1, 2, \cdots, n, \\ |\delta_{ii}| &\leq |\varepsilon_{ii}| + |\varepsilon'_{ii}| + |\varepsilon_{ii}\varepsilon'_{ii}| \leq 2 \cdot 2^{-i} + 2^{-2i}. \end{aligned} \quad (2.2.48)$$

由(2.2.46), (2.2.48)可知

$$\begin{aligned} |\delta U| &\leq (1.06) 2^{-1} \\ &\begin{pmatrix} |u_{11}| & & & \\ 1|u_{21}| & 2|u_{22}| & & \\ 2|u_{31}| & 2|u_{32}| & 2|u_{33}| & \\ \cdots & \cdots & \cdots & \cdots \\ (n-1)|u_{n1}| & (n-1)|u_{n2}| & \cdots & 2|u_{nn-1}| & 2|u_{nn}| \end{pmatrix} \end{aligned} \quad (2.2.49)$$

估计 δU 的范数 $\|\cdot\|_\infty$ 得到

$$\begin{aligned} \|\delta U\| &\leq (1.06) 2^{-1} \left(\sum_{i=1}^n i \right) \max_{i,k} |u_{ik}| \\ &= \frac{n^2 + n}{2} (1.06) \cdot 2^{-1} \max_{i,k} |u_{ik}|. \end{aligned} \quad (2.2.50)$$

定理2.2.2 系数矩阵为三角阵的线性方程组 $Ux = b$ 的计算解是方程组 $(U + \delta U)x = b$ 的精确解, 其中 δU 满足(2.2.49)和(2.2.50).

将这个定理应用到求解三角方程组 $\tilde{L}y = b$ 和 $\tilde{R}x = y$, 最后得到的计算解应满足

$$(\tilde{L} + \delta\tilde{L})(\tilde{R} + \delta\tilde{R})x = b$$

或 $(\tilde{L}\tilde{R} + \delta\tilde{L}\cdot\tilde{R} + \tilde{L}\cdot\delta\tilde{R} + \delta\tilde{L}\cdot\delta\tilde{R})x = b.$

将 $\tilde{L}\tilde{R} = A + E$ 代入上式, 得

$$(A + E + \delta\tilde{L}\cdot\tilde{R} + \tilde{L}\cdot\delta\tilde{R} + \delta\tilde{L}\cdot\delta\tilde{R})x = b, \quad (2.2.51)$$

由(2.2.36), (2.2.37)和(2.2.50)知

$$\left\{ \begin{aligned} \|\tilde{L}\|_{\infty} &\leq n, \\ \|\tilde{R}\|_{\infty} &\leq n\rho\|A\|_{\infty}, \\ \|\delta\tilde{L}\|_{\infty} &\leq \frac{n^2+n}{2}(1.06)\cdot 2^{-t}\max_{i,k}|l_{i,k}| \leq \frac{n^2+n}{2}(1.06)2^{-t}, \\ \|\delta\tilde{R}\|_{\infty} &\leq \frac{n^2+n}{2}(1.06)\cdot 2^{-t}\max_{i,h}|r_{i,h}| \\ &\leq \frac{n^2+n}{2}(1.06)\rho\|A\|_{\infty}2^{-t}. \end{aligned} \right. \quad (2.2.52)$$

现今

$$\delta A = E + \delta\tilde{L}\cdot\tilde{R} + \tilde{L}\cdot\delta\tilde{R} + \delta\tilde{L}\cdot\delta\tilde{R}, \quad (2.2.53)$$

则 $\|\delta A\|_{\infty} \leq \|E\|_{\infty} + \|\delta\tilde{L}\|_{\infty}\|\tilde{R}\|_{\infty} + \|\tilde{L}\|_{\infty}\|\delta\tilde{R}\|_{\infty} + \|\delta\tilde{L}\|_{\infty}\|\delta\tilde{R}\|_{\infty}.$

由(2.2.40), (2.2.52)可得

$$\|\delta A\|_{\infty} \leq 1.06n^2\rho\|A\|_{\infty}\cdot 2^{-t}\left(3+n+\frac{(n+1)^2}{2}2^{-t}\right).$$

又因为, 在实际计算中, 常有 $(n+1)^2\cdot 2^{-t}/2 \ll 1$, 故有

$$\|\delta A\|_{\infty} \leq 1.06(n^3+4n^2)\rho\cdot 2^{-t}\|A\|_{\infty} \quad (2.2.54)$$

定理2.2.3 用列主元消去法求出的线性方程组 $Ax = b$ 的计算解 x 是方程组

$$(A + \delta A)x = b$$

的精确解. 其中误差矩阵 δA 由(2.2.53)给出, 且满足估计式(2.2.54).

从定理 2.2.3 的证明过程中, 我们可以看到, 由 (2.2.54) 所表示的误差矩阵的范数的上界估计式还是较保守的。如果令 $\bar{A} = A + \delta A$, 则定理 2.2.3 的结果表明 $\phi^*(A, b) = \phi(\bar{A}, b)$, 其中 $\phi^*(A, b)$ 表示 $Ax = b$ 的计算解, $\phi(\bar{A}, b)$ 表示 $\bar{A}x = b$ 的精确解, 而且有 $\|\bar{A} - A\|_{\infty} \leq 1.06(n^3 + 4n^2)\rho \cdot 2^{-t}\|A\|_{\infty}$ 。只要 $1.06(n^3 + 4n^2)\rho\|A\|_{\infty}$ 有一个不大的常数上界, 即 $|\delta A|$ 的元素很小, 那末, 计算解 ϕ^* (按列的部份主元素法) 就是稳定的。在 $1.06(n^3 + 4n^2)\rho\|A\|_{\infty}$ 中, 需要进一步分析的是 ρ 的上界的估计问题。

由于

$$\begin{aligned} |a_{ij}^{(k)}| &\leq |a_{ij}^{(k-1)} - m_{ik}a_{kj}^{(k-1)}| \leq |a_{ij}^{(k-1)}| + |a_{kj}^{(k-1)}| \\ &\leq 2\max_{i,j} |a_{ij}^{(k-1)}|, \end{aligned}$$

从而有

$$\begin{aligned} \max_{i,j} |a_{ij}^{(k)}| &\leq 2\max_{i,j} |a_{ij}^{(k-1)}| \leq 2^2\max_{i,j} |a_{ij}^{(k-2)}| \leq \dots \\ &\leq 2^{k-1}\max_{i,j} |a_{ij}| \leq 2^{k-1}\|A\|_{\infty}, \end{aligned}$$

因此

$$\rho \leq m \cdot 2^{n-1}, \quad m = \max_{i,j} |a_{ij}| / \|A\|_{\infty} \quad (\text{常数}).$$

从而 ρ 有一个与矩阵阶数 n 有关的上界, 其中 2^{n-1} 是一个变化相当快的增长因子。但是在实践中很少达到这个界, 因此在实践中解线性方程组的列主元素法被认为是稳定的。

2.3 解的精度

定理 2.2.3 说明了非奇异线性方程组 $Ax = b$ 的计算解 x 正好是方程组

$$(A + \delta A)x = b$$

的精确解, 并且给出了误差矩阵 δA 的界为

$$\|\delta A\|_{\infty} < 1.06(n^3 + 4n^2)\rho \cdot 2^{-t}\|A\|_{\infty}.$$

如果令 $f(n) = 1.06(n^3 + 4n^2)\rho$, 这是一个依赖于具体计算细节的量, 一般说来, 其值不会很大, 将 $f(n)$ 代入上式得到

$$\|\delta A\|_{\infty} < f(n) \|A\|_{\infty} \cdot 2^{-t}, \quad (2.2.55)$$

由定理 2.1.2 的推论知道, 如果设 \bar{x} 是 $Ax = b$ 的精确解, 且 $\|A^{-1}\|_{\infty} \|\delta A\|_{\infty} < 1$, 则有

$$\frac{\|\bar{x} - x\|_{\infty}}{\|\bar{x}\|_{\infty}} \leq \frac{\|A^{-1}\|_{\infty} \|\delta A\|_{\infty}}{1 - \|A^{-1}\|_{\infty} \|\delta A\|_{\infty}}.$$

因此, 当

$$\begin{aligned} \|A^{-1}\|_{\infty} \|\delta A\|_{\infty} &< f(n) \|A^{-1}\|_{\infty} \|A\|_{\infty} \cdot 2^{-t} \\ &= f(n) \operatorname{cond}(A) 2^{-t} < 1 \end{aligned} \quad (2.2.56)$$

时, 就有

$$\frac{\|\bar{x} - x\|_{\infty}}{\|\bar{x}\|_{\infty}} \leq \frac{f(n) \operatorname{cond}(A) 2^{-t}}{1 - f(n) \operatorname{cond}(A) 2^{-t}}. \quad (2.2.57)$$

如果 $f(n)$ 和 $\operatorname{cond}(A)$ 均不太大, 并且令 $f(n) \operatorname{cond}(A) \leq 2^p$, 则

$$\frac{\|\bar{x} - x\|_{\infty}}{\|\bar{x}\|_{\infty}} \leq \frac{2^{-t+p}}{1 - 2^{-t+p}}, \quad (2.2.58)$$

从而我们可以期望用 t 位精度求得的计算解, 能达到大约有 $k = t - p$ 位有效数字. 而且对 $\|x\|_{\infty}$ 有如下的估计式

$$\left(1 - \frac{2^{-k}}{1 - 2^{-k}}\right) \|\bar{x}\|_{\infty} \leq \|x\|_{\infty} \leq \frac{\|\bar{x}\|_{\infty}}{1 - 2^{-k}}. \quad (2.2.59)$$

但是必须注意, 如果条件 (2.2.56) 不满足时, 我们就不能肯定 x 有任何精度, 而且 $(A + \delta A)$ 甚至可能是奇异的. 由于对 $\|\delta A\|_{\infty}$ 的上界的估计是非常保守的, 因此, 利用 (2.2.57) 来估计解的精确度时难免会产生过估现象, 有时候这种过估还是很严重的.

由于 $f(n)$ 和 $\text{cond}(A)$ 均有可能进行估计, 从而在算得 x 以前就可以定出 x 的相对误差界, 即先验误差界. (2.2.57) 中的误差界即属于先验误差界. 现在再来分析依赖于近似解 x 本身的误差界, 即通常所称的事后误差界. 观察事后误差界的一个很自然的方法是把它代入方程组 $Ax = b$, 然后再看残向量

$$r = b - Ax \quad (2.2.60)$$

范数的大小. 对于 x , 如果 $r = 0$, 则 x 就是 $Ax = b$ 的精确解. 然而, 是否 $\|r\|$ 较小时, 近似解 x 的精确度就一定高呢? 我们可以从下面的定理得出应有的结论.

定理 2.2.4 设 A 是非奇异阵且 $Ax = b \neq 0$, x 是已知的近似解, \bar{x} 为精确解, 则

$$\frac{\|x - \bar{x}\|}{\|\bar{x}\|} \leq \text{cond}(A) \frac{\|r\|}{\|b\|}. \quad (2.2.61)$$

证明 因为 $r = b - Ax$, 且 A 非奇异, 所以

$$A^{-1}r = A^{-1}b - x = \bar{x} - x,$$

于是

$$\|\bar{x} - x\| = \|A^{-1}r\| \leq \|A^{-1}\| \|r\|, \quad (2.2.62)$$

又
故

$$\|b\| = \|A\bar{x}\| \leq \|A\| \|\bar{x}\|$$

$$\|\bar{x}\| \geq \|b\| / \|A\|. \quad (2.2.63)$$

(2.2.62) 除以 (2.2.63) 得

$$\frac{\|\bar{x} - x\|}{\|\bar{x}\|} \leq \frac{\|A^{-1}\| \|A\| \|r\|}{\|b\|} = \text{cond}(A) \frac{\|r\|}{\|b\|}.$$

这个结果表明, 解的精度不仅依赖于残向量的模 $\|r\|$ 的大小, 而且还依赖于矩阵的条件数. 如果 A 是病态的, 即使残向量的

模 $\|r\|$ 很小，也不能保证近似解的精度很高，而且对于精度高的近似解也可能有大的 $\|r\|$ 。

例2.2.2 方程组

$$\begin{bmatrix} 1.000 & 1.001 \\ 1.000 & 1.000 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2.001 \\ 2.000 \end{bmatrix}$$

的精确解显然是 $\bar{x} = (1, 1)^T$ 。向量 $x = (2, 0)^T$ 并不接近 \bar{x} ，而残向量

$$r = \begin{bmatrix} 2.001 \\ 2.000 \end{bmatrix} - \begin{bmatrix} 1.000 & 1.001 \\ 1.000 & 1.000 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.001 \\ 0.000 \end{bmatrix}$$

却很小。此例说明残向量很小，但近似解的精度并不一定高。

方程组

$$\begin{bmatrix} 1.000 & 1.001 \\ 1.000 & 1.000 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

的精确解显然是 $\bar{x} = [-1000, 1000]^T$ ，但是向量 $(-1001, 1000)^T$ 虽然较接近于 \bar{x} ，相应的残向量

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 1.000 & 1.001 \\ 1.000 & 1.000 \end{bmatrix} \begin{bmatrix} -1001 \\ 1000 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = r$$

却比 b 还大。此例说明近似解的精度虽然较高，但是残向量并不一定很小。

什么原因引起方程组病态以及怎样度量一个方程组是否病态呢？这一直是人们关心的一个困难的课题。Turing在1948年第一个使用了“条件数”这个术语。接着，不少数学工作者相继提出用条件数来度量相应矩阵的病态的程度，但是对条件数的大小的估计还缺少有效的方法。Noble 1969年所给出的检查病态条件的各种方法的图表，使用起来仍不够理想。随着电子计算机的发展，解病态问题也变得愈来愈突出，近年来，除解

病态方程组的直接法和迭代法仍在继续发展外，还出现了预处理法，这是一种较有发展前途的方法，但是目前还都有一些局限性。因此，判别病态问题和解病态问题都是目前计算数学的重要课题。

关于算法的误差分析，经 Wilkinson 的努力，自 1957 年取得了比较大的进展后，直到目前向后误差分析还是分析舍入误差的最完整最系统的方法。近年来，它的概念更加系统，分析过程趋于简化，结果也更为精确而实用。

第二章 习 题

2.1 证明方程组

$$\begin{cases} x + y = 2, \\ \frac{1}{6}x + \frac{1}{6}y = \frac{1}{3}, \\ x + 2y = 3 \end{cases}$$

有唯一解，但是它不适定。

2.2 设 $A \in R^{n \times n}$ 为 n 阶矩阵， $x = (x_1, x_2, \dots, x_n)^T$ ，其中 a_{ij} 和 x_i 均为规格化浮点数， k 为某实数，求 $fl(kA)$ ， $fl(Ax)$ 。

2.3 $A, B \in R^{n \times n}$ 其中 a_{ij} 和 b_{ij} 均为规格化浮点数，求 $fl(A \cdot B)$ ， $fl(A + B)$ 。

2.4 设 x, y, z 都是具有 t 个数字的规格化浮点数，用例子说明

$$\begin{aligned} fl(fl(x+y) + z) &\neq fl(x + fl(y+z)), \\ fl(fl(xy)z) &= fl(x fl(yz)) \end{aligned}$$

是否一定成立？为什么？

2.5 设 $fl(x_1 + \dots + x_k) \equiv fl(fl(x_1 + \dots + x_{k-1}) + x_k)$

$$(k = 2, \dots, n),$$

证明 $fl(a_1 + a_2 + \dots + a_n) = a_1(1 + \rho_1)^{n-1} +$

$$a_2(1 + \rho_2)^{n-1} + a_3(1 + \rho_3)^{n-2} + \dots + a_n(1 + \rho_n),$$

其中 $|\rho_i| \leq u 10^{-1}$, u 是与 1 同阶的数。

2.6 设 $x = 39.25$, $y = 0.23$, 试将 x , y 分别化成 $t = 6$ 的规格化浮点数的机器表示形式 $a = 2^b a_1$ 和 $b = 2^{b_1} a_2$, 并求 $fl(a + b)$, $fl\left(\frac{a}{b}\right)$ 。

2.7 设二次方程

$$x^2 - 6.433x + 0.009474 = 0,$$

试用求根公式计算该方程的最小根的计算值。

2.8 举出“病态”和不稳定算法的数值例子各一个。

2.9 如果系数矩阵 A 是三对角阵, 使用部分主元素法求解 $Ax = b$ 时, 估计 $\rho = \max_{i,j,k} |a_{ij}^{(k)}| / \|A\|_\infty$ 的上界。

2.10 设 A 为实对称矩阵, R 为正交阵, 求 $fl(R^T A R)$, 并给出误差矩阵在某种范数下的上界。

2.11 证明定理 2.1.1 即证明关系式 $k^{(1)}(A) \leq \text{cond}(A) \leq \sqrt{n} 2^{n-1} k^{(1)}(A)$ 。

2.12 作镜像变换的误差分析。

参 考 书

- [1] J. H. Wilkinson, "The Algebraic eigenvalue problems", Oxford Univ. press, London and New York, 1965.
- [2] Stewart. G, "Introduction to Motrix Computation", Academic press, New York 1973 (中译本, 矩阵计算引论, 王国荣等译, 上海科技出版社)。
- [3] 曹志浩等编, 矩阵计算和方程求根, 1979. 人民教育出版社。
- [4] G. M. phillips, P. J. Taylor, "Theory and Applications of Numerical Analysis", Academic press, 1973 (中译本, 数值分析的理论及其应用, 熊西文等译, 上海科技出版社)。

第三章 特征值计算

§1 引言

1.1 相似变换

用矩阵变换方法求矩阵的特征值，其基本思想是根据相似矩阵必有相同的特征谱这一性质，把矩阵化为形式简单的矩阵，使新矩阵的特征值便于计算。上三角阵是形式简单的矩阵，且其对角元就是它的特征值。

如果只限于用实数进行运算，那末，有复特征值的实矩阵在任何相似变换下都不能化为上三角阵。比上三角阵范围更广且特征值又易于计算的一类形式简单的矩阵是形如(3.1.1)的被称为拟上三角阵的矩阵：

$$\begin{pmatrix} \times & \times & & & \\ \times & \times & & & \\ & & \begin{matrix} \times & \times \\ \times & \times \end{matrix} & & \\ & & & \ddots & \\ & & & \times & \\ & & & & \times \end{pmatrix} \quad (3.1.1)$$

这种矩阵的下三角部分，除次对角线上有非零元素外，其余全为零，而次对角线上的两个非零元素之间必有一些零元，这是一个对角线上由一阶或二阶块组成的阵。而矩阵特征值就是这些块阵的特征值组成的。于是，求给定矩阵的特征值问题，就

可归结为：把给定矩阵用一系列相似变换化为上三角阵或拟上三角阵。

这个过程，通常是不能在有限步内结束的。例如，形如

$$F = \begin{pmatrix} -p_1 & -p_2 & \cdots & -p_n \\ & 1 & 0 & \\ & & 1 & \ddots \\ & & & \ddots & 1 & 0 \end{pmatrix} \quad (3.1.2)$$

的所谓 Frobenious 阵，其特征多项式为

$$f(\lambda) = \lambda^n + p_1 \lambda^{n-1} + \cdots + p_{n-1} \lambda + p_n,$$

当 $n > 5$ 时，它的零点，即 F 的特征值，一般是不能用 $p_i (i = 1, 2, \cdots, n)$ 经有限次算术运算求出的。由此可知，使 $P^{-1}FP$ 化为 (3.1.1) 的 P 也不能由 $p_i (i = 1, 2, \cdots, n)$ 经有限次运算求得。

使任意实矩阵 A 化为拟三角阵的相似变换矩阵是否存在与这个变换矩阵是否可经由有限步运算得到，这是两个不同的问题，后者已予说明，前者可用下列定理来回答。

定理 3.1.1 设 A 是实矩阵，则存在非奇异的实矩阵 P ，使 $P^{-1}AP$ 成为拟上三角阵。

证明 用归纳法，当 $n = 1$ 时，定理显然成立。设当 $n \leq k$ 时定理也成立。则当 $n = k + 1$ 时，任取 A 的一个特征值 λ ，当 λ 是实数时，可设其相应于 λ 的特征向量为 x ，且 x 也是实的，这时，可以找到一个实的初等 Hermite 矩阵 H_1 ，此时 H_1 是实的正交阵，使

$$H_1 x = \rho e_1, \quad |\rho| = \|x\| \neq 0 \quad (3.1.3)$$

由于

$$Ax = \lambda x$$

两端用 H_1 左乘, 并注意 $\mathbf{x} = PH_1^{-1}\mathbf{e}_1$, 得

$$H_1AH_1^{-1}\mathbf{e}_1 = \lambda\mathbf{e}_1, \quad (3.1.4)$$

可见 $H_1AH_1^{-1}$ 具有下列形状.

$$H_1AH_1^{-1} = \left(\begin{array}{c|c} \lambda & \mathbf{b}_1^T \\ \hline \mathbf{0} & A_{n-1} \end{array} \right). \quad (3.1.5)$$

据归纳法假设, A_{n-1} 是阶数 $\leq k$ 的阵, 存在正交阵 H_2 使 $H_2A_{n-1}H_2^{-1}$ 具有 (3.1.1) 的形状, 取

$$P = \left(\begin{array}{c|c} 1 & \mathbf{0} \\ \hline \mathbf{0} & H_2 \end{array} \right) H_1, \quad (3.1.6)$$

则 PAP^{-1} 亦具有 (3.1.1) 的形状. 且 P 是非奇异阵.

当 λ 是复数时, 令 $\lambda = \alpha + i\beta$, $\mathbf{x} = \mathbf{u} + i\mathbf{v}$, 而 $\alpha, \beta, \mathbf{u}, \mathbf{v}$ 都是实的, 且 $\beta \neq 0$, 由 $A\mathbf{x} = \lambda\mathbf{x}$, 得

$$\begin{cases} A\mathbf{u} = \alpha\mathbf{u} - \beta\mathbf{v}, \\ A\mathbf{v} = \beta\mathbf{u} + \alpha\mathbf{v}, \end{cases} \quad (3.1.7)$$

可以证明 \mathbf{u}, \mathbf{v} 是线性无关的. 不然, 设 p, q 不全为零, 而 $p\mathbf{u} + q\mathbf{v} = \mathbf{0}$, 则由 (3.1.7), 可得

$$\beta(-p\mathbf{v} + q\mathbf{u}) = \mathbf{0},$$

从而有

$$p\mathbf{u} + q\mathbf{v} = \mathbf{0} \quad \text{及} \quad q\mathbf{u} - p\mathbf{v} = \mathbf{0}.$$

由于 $p^2 + q^2 \neq 0$, 可得 $\mathbf{u} = \mathbf{v} = \mathbf{0}$, 矛盾. 故 \mathbf{u}, \mathbf{v} 线性无关. 这时, 可以找到初等正交阵 H_1 , 使

$$H_1(\mathbf{u}, \mathbf{v}) = \begin{pmatrix} r_{11} & r_{12} \\ 0 & r_{22} \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{pmatrix}, \quad r_{11}r_{22} \neq 0, \quad (3.1.8)$$

据(3.1.7), 有

$$A(u, v) = (u, v) \begin{pmatrix} \alpha & -\beta \\ \beta & \alpha \end{pmatrix} \quad (3.1.9)$$

及

$$H_1 A H_1^{-1} \begin{pmatrix} r_{11} & r_{12} \\ 0 & r_{22} \\ \vdots & \vdots \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} r_{11} & r_{22} \\ 0 & r_{22} \\ \vdots & \vdots \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \alpha & -\beta \\ \beta & \alpha \end{pmatrix}. \quad (3.1.10)$$

比较上式两端, 知 $H_1 A H_1^{-1}$ 具有下列形状

$$H_1 A H_1^{-1} = \begin{pmatrix} A_2 & B_2 \\ 0 & A_{n-2} \end{pmatrix}. \quad (3.1.11)$$

这里

$$A_2 = \begin{pmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{pmatrix} \begin{pmatrix} \alpha & -\beta \\ \beta & \alpha \end{pmatrix} \begin{pmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{pmatrix}^{-1},$$

可见 A_2 的特征值仍为 $\alpha \pm i\beta$, 由于 A_{n-2} 的阶是 $k-1$, 据归纳法假设, 仿前面的讨论, 知定理对 $n=k+1$ 也成立.

定理 3.1.1 的证明, 将以 A 的特征值及特征向量为前提, 并不能看作是使 $P^{-1}AP$ 成为拟三角阵的 P 的存在的构造性证明, 但定理本身却为 P 的构造性算法提供了理论根据.

1.2 LR 算法

设实矩阵 A 有三角分解

$$A = LR, \quad (3.1.12)$$

其中 L 是单位下三角阵, R 是上三角阵, L 显然有逆, 且 $L^{-1}AL$ 是 A 的一个相似矩阵, 记

$$A_2 = L^{-1}A_1L, \quad A_1 = A, \quad (3.1.13)$$

代入 $A = LR$, 则

$$A_2 = RL, \quad (3.1.14)$$

它不过是 A 的三角分解中两个矩阵按相反的次序的乘积, A_2 是新矩阵, 用它取代 $A_1 = A$ 的地位, 并重复刚才的过程, 又可得到 A_3 , 仿此, 便可得到与 A 相似的矩阵序列:

$$\begin{aligned} A_1 &= A, \quad A_s = L_s R_s, \quad A_{s+1} = R_s L_s, \\ s &= 1, 2, \dots \end{aligned} \quad (3.1.15)$$

这里每个 L_s 都是单位下三角阵, 每个 R_s 都是上三角阵。

注意, 并非每个矩阵 A , 甚至非奇异的矩阵 A 都有 LR 分解, 例如

$$A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

是非奇异的, 但却没有 LR 分解。

当 A 非奇异时, 据第一章, A 有 LR 分解的充要条件应是 A 的各阶首主子式不为零, 在此条件下分解将是唯一的。本节在讨论 A 的 LR 分解时, 我们总是假定 A 及其 LR 分解序列中的每一矩阵 A_s 的各阶首主子式均不为零。

设 $\{A_s\}$, $s = 1, 2, \dots$ 是 A 的 LR 分解序列, 则可证:

$$A_1 = L_1 L_2 \cdots L_{s-1} A_s L_s^{-1} \cdots L_2^{-1} L_1^{-1}, \quad (3.1.16)$$

$$A_1 L_1 \cdots L_{s-1} = L_1 \cdots L_{s-1} A_s. \quad (3.1.17)$$

记

$$T_s = L_1 L_2 \cdots L_s, \quad U_s = R_s \cdots R_1 \quad (3.1.18)$$

$$\text{则有} \quad A_1^s = T_s U_s, \quad (3.1.19)$$

这表明 T_s, U_s 是 A_1^s 的 LR 分解。

1.3 LR 算法的收敛性

由 $A_1^s = T_s U_s$ 及 A_1 的非奇异性, 知 A_1^s 的三角分解应是唯一的。记 $t_{ij}^{(s)}$ 为 T_s 中的 (i, j) 元素, 其中 $t_{ii}^{(s)} = 1$, $t_{ij}^{(s)} = 0$,

$(j > i)$, 令 $\Delta_i^{(s)}$ 表示 A_1^s 的 i 阶首主子式, $\Delta_i^{(s)}$ 表示 $\Delta_i^{(s)}$ 中把第 j 行的元素换成 A_1^s 的第 i 行 $(i > j)$ 上相应的元素后所成的

j 阶行列式, 则有

$$t_{ij}^{(s)} = \Delta_i^{(s)} / \Delta_j^{(s)} \quad (i > j). \quad (3.1.20)$$

下面考查 A 在什么条件下, 其 LR 序列是收敛的.

设 A 有依模不相等的特征值.

$$|\lambda_1| > |\lambda_2| > \cdots > |\lambda_n|, \quad (3.1.21)$$

对应的特征向量为 x_1, x_2, \cdots, x_n , 记

$$X = (x_1, \cdots, x_n) = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{pmatrix}, \quad (3.1.22)$$

记 X 的逆阵为 Y :

$$X^{-1} = Y(y_1, \cdots, y_n) = \begin{pmatrix} y_{11} & y_{12} & \cdots & y_{1n} \\ y_{21} & y_{22} & \cdots & y_{2n} \\ \vdots & \vdots & & \vdots \\ y_{n1} & y_{n2} & \vdots & y_{nn} \end{pmatrix}, \quad (3.1.23)$$

于是有

$$A = X \operatorname{diag}(\lambda_1, \cdots, \lambda_n) Y,$$

$$A^s = X \operatorname{diag}(\lambda_1^s, \cdots, \lambda_n^s) Y,$$

及

$$(A^s)_{ij} = \sum_{k=1}^n x_{ik} y_{kj} \lambda_k^s, \quad (3.1.24)$$

A^s 的 p 阶首主子式为

$$\begin{aligned} \Delta_p^{(s)} = & \begin{vmatrix} \sum_{k=1}^n x_{1k} y_{k1} \lambda_k^s & \sum_{k=1}^n x_{1k} y_{k2} \lambda_k^s & \cdots & \sum_{k=1}^n x_{1k} y_{kp} \lambda_k^s \\ \sum_{k=1}^n x_{2k} y_{k1} \lambda_k^s & \sum_{k=1}^n x_{2k} y_{k2} \lambda_k^s & \cdots & \sum_{k=1}^n x_{2k} y_{kp} \lambda_k^s \\ \vdots & \vdots & & \vdots \\ \sum_{k=1}^n x_{pk} y_{k1} \lambda_k^s & \sum_{k=1}^n x_{pk} y_{k2} \lambda_k^s & \cdots & \sum_{k=1}^n x_{pk} y_{kp} \lambda_k^s \end{vmatrix} \\ & = X_p \bar{Y}_p (\lambda_1, \cdots, \lambda_p)^s + \cdots \end{aligned} \quad (3.1.25)$$

这里 X_p 、 Y_p 分别表示 X 、 Y 的 p 阶首主子式，上式中没有明显写出的项是形如

$$X_{k_1 \dots k_p} \bar{Y}_{k_1 \dots k_p} (\lambda_{k_1}, \dots, \lambda_{k_p})^s \\ | \leq k_1 < k_2 < \dots < k_p \leq n$$

的项，而 $X_{k_1 \dots k_p}$ 表示 X 的一个 p 阶子式，由 X 的 1 至 p 行和第 $k_1 \dots k_p$ 诸列上的元素组成， $Y_{k_1 \dots k_p}$ 是 Y 的一个 p 阶子式，由 \bar{Y} 的 1 至 p 列及第 k_1, k_2, \dots, k_p 各行上的元素组成， k_1, k_2, \dots, k_p 是 $1, 2, \dots, n$ 的任意一个元素个数为 p 的子集，当 k_1, k_2, \dots, k_p 分别是 $1, \dots, p$ 时，简记为 X_p 和 Y_p 。当 $X_p Y_p \neq 0$ 时，(3.1.25) 右端中的第一项将是 $\Delta_p^{(s)}$ 的主要部分。

记 $\Delta_{i,p}^{(s)}$ 和 $X_p^{(i)}$ 分别为将 $\Delta_p^{(s)}$ 和 X_p 中的第 p 行换为 A^s 和 X 的 i 行后所成的行列式，则

$$\Delta_{i,p}^{(s)} = X_p^{(i)} Y_p (\lambda_1, \dots, \lambda_p)^s + \dots, \quad (3.1.26)$$

于是当 $X_p \bar{Y}_p \neq 0$ 时，有

$$t_{i,p}^{(s)} = \frac{X_p^{(i)} Y_p (\lambda_1, \dots, \lambda_p)^s + \dots}{X_p \bar{Y}_p (\lambda_1, \dots, \lambda_p)^s + \dots} \\ = \frac{X_p^{(i)}}{X_p} + O\left(\frac{\lambda_{p+1}}{\lambda_p}\right)^s, \quad (3.1.27)$$

当 $s \rightarrow \infty$ 时， $t_{i,p}^{(s)}$ 有极限，而当 $X_p = 0$ ， $Y_p \neq 0$ ， $X_p^{(i)} \neq 0$ 时， $t_{i,p}^{(s)}$ 将是发散的。

据此，乃有下列关于矩阵的 LR 分解序列收敛的定理。

定理 3.1.2 设 A 有 LR 分解， A 的特征值依模不等，并且 A 的特征向量所成的矩阵以及该矩阵的逆的各阶首主子式不为零，则 A 的 LR 分解序列必收敛于某一上三角阵。而极限矩阵的对角元就是矩阵 A 的特征值。

证明 定理的条件保证了 $t_{i,p}^{(s)}$ ($i > p$) 的收敛性，但 (3.1.27)

右端的极限恰是 X 的 LR 分解中的下三角阵的 (i, p) 元素, 即

$$T_i \rightarrow T \quad (3.1.28)$$

这里 T 是 X 的 LR 分解 $X = TS$ 中的单位下三角阵

$$\begin{aligned} A_i &= T_{i-1}^{-1} A T_{i-1} \rightarrow T^{-1} A T \\ &= S \text{diag}(\lambda_i) S^{-1}. \end{aligned} \quad (3.1.29)$$

上式右端是以 $\lambda_1, \lambda_2, \dots, \lambda_n$ 为对角的上三角阵。

当 $A = A_1$ 是对称正定阵时, 由于 X 可取为正交阵, 有 $Y = X^{-1} = X^T$, 可得

$$X_{i_1 \dots i_p} = Y_{i_1 \dots i_p}$$

及

$$t_{i,p}^{(s)} = \frac{\sum X_{i_1 \dots i_p}^{(i)} Y_{i_1 \dots i_p} (\lambda_{i_1}, \dots, \lambda_{i_p})^s}{\sum (X_{i_1 \dots i_p})^2 (\lambda_{i_1}, \dots, \lambda_{i_p})^s} \quad (3.1.30)$$

这里 $X_{i_1 \dots i_p}^{(i)}$ 表示 $X_{i_1 \dots i_p}$ 中把第 p 行的元素换为 X 的第 i 行 ($i > p$) 的元素后所成的行列式。将 (3.1.30) 两端取极限, 可见, 当 $s \rightarrow \infty$ 时, $t_{i,p}^{(s)}$ 不会遇到发散的情况, 记

$$T_i \rightarrow T_\infty \quad (3.1.31)$$

此时 T_∞ 未必是 T , 但仍是单位下三角阵, 而由 (3.1.31)

$$L_i = T_{i-1}^{-1} T_i \rightarrow T_\infty^{-1} T_\infty = I$$

但

$$\begin{aligned} A_i &= T_{i-1}^{-1} A T_{i-1} \rightarrow T_\infty^{-1} A_1 T_\infty \\ R_i &= L_i^{-1} A_i \rightarrow T_\infty^{-1} A_1 T_\infty \end{aligned} \quad (3.1.32)$$

可见 A_i 与 R_i 有同样的极限, 由于 R_i 是上三角阵, 故其极限也应是上三角阵, 并且该极限阵的对角元应是原矩阵 A 的特征值。

§ 2 QR 算法及收敛性

2.1 QR 算法的基本收敛性

由于矩阵的 LR 分解中, 其单位下三角阵中的元素不能保证依模小于 1, 甚至无法估计出这些元素模的上界, 因此, LR

算法的稳定性便没有保证, 实践经验证明, 在这种算法下, 计算特征值时精度的损失, 通常并不是偶然的现象。1961年以来, **Francis** 对此法进行了改进, 用酉变换作工具, 用 **QR** 分解代替 **LR** 分解而发展起来的 **QR** 方法, 在大多数情况下被证明是求一般代数特征值问题的最有效的方法。

令 $A = A_1$ 是非奇异实阵, 则 A 可表示为

$$A_1 = Q_1 R_1, \quad (3.2.1)$$

其中 Q_1 是正交阵, R_1 是上三角阵, 当规定 R_1 的对角元是非负的元素时, 此种分解将是唯一的。令

$$A_2 = R_1 Q_1, \quad (3.2.2)$$

易证 A_2 相似于 A_1 , 一般有

$$A_1 = A, \quad A_s = Q_s R_s, \quad A_{s+1} = R_s Q_s, \quad (3.2.3)$$

$$s = 1, 2, \dots$$

这里每个 Q_s 仍是正交阵, R_s 仍是对角元非负的上三角阵。不难证明

$$A_{s+1} = Q_s^H \cdots Q_1^H A_1 Q_1 \cdots Q_s, \quad (3.2.4)$$

$$Q_1 \cdots Q_s A_{s+1} = A_1 Q_1 \cdots Q_s. \quad (3.2.5)$$

令

$$Q_1 \cdots Q_s = P_s, \quad P_s \cdots P_1 = U_s,$$

则

$$A_1^s = P_s U_s. \quad (3.2.6)$$

在讨论 A 的 **QR** 分解序列的收敛性之前, 先考察一个例子。

例 3.2.1 令

$$A = \begin{pmatrix} e^{i\theta} & 1 \\ & e^{-i\theta} \end{pmatrix},$$

则其 **QR** 分解序列可表为

$$A_s = \begin{pmatrix} e^{i\theta} & e^{-2(s-1)i\theta} \\ & e^{-i\theta} \end{pmatrix}.$$

这里，由于 $(1, 2)$ 元素在 $s \rightarrow \infty$ 时无极限，所以 A_s 在通常的意义下是不收敛的。然而，对于 A_s 的下三角阵收敛于对角阵的情况，不管 A_s 的严格上三角阵有无极限，只要相应的元素保持有界，且把相应的对角阵当作 A_s 的极限，则 A_s 在实质上收敛于此对角阵。在本章后面的讨论中，当谈到矩阵列的极限时，我们都把它当作依实质收敛来理解。

下面我们来讨论 A 的 QR 分解序列 $\{A_s\}$ 收敛的充分条件。

定理 3.2.1 设 A 非奇异且其特征值依模不等：

$$|\lambda_1| > |\lambda_2| > \cdots > |\lambda_n| > 0 \quad (3.2.7)$$

设其特征向量所成的矩阵为 X ，当 X 的逆 Y 有 LR 分解，即

$$X^{-1} = Y = L_y U_y \quad (3.2.8)$$

其 $\{A_s\}$ 将在实质上收敛于依次以 $\lambda_1, \lambda_2, \dots, \lambda_n$ 为对角元的上三角阵。

证明 由

$$\begin{aligned} A &= X \text{diag}(\lambda_1, \dots, \lambda_n) Y = X \Lambda Y, \\ A^s &= X \Lambda^s Y, \quad \Lambda = \text{diag}(\lambda_1 \cdots \lambda_n), \end{aligned} \quad (3.2.9)$$

代入 X 的 QR 分解及 Y 的 LR 分解：

$$X = Q_x R_x, \quad Y = L_y U_y, \quad (3.2.10)$$

则

$$A^s = Q_x R_x \Lambda^s L_y U_y = Q_x R_x (\Lambda^s L_y \Lambda^{-s}) (\Lambda^s U_y),$$

易见 $\Lambda^s L_y \Lambda^{-s}$ 仍是单位下三角阵，其 (i, j) 元素 ($j < i$)

$$l_{ij} (\lambda_i / \lambda_j)^s, \quad (j < i) \quad (3.2.11)$$

是趋于零的，这里 l_{ij} 是 L_y 的 (i, j) 元素。记

$$\Lambda^s L_y \Lambda^{-s} = I + E_s, \quad E_s \rightarrow 0, \quad (3.2.12)$$

则

$$\begin{aligned} A^s &= Q_x R_x (I + E_s) \Lambda^s U_y \\ &= Q_x (I + R_x E_s R_x^{-1}) R_x \Lambda^s U_y, \end{aligned}$$

记

$$D_1 = \begin{pmatrix} e^{i\theta_1} & & \\ & e^{i\theta_2} & \\ & & \ddots \\ & & & e^{i\theta_n} \end{pmatrix}, \quad \tilde{\Lambda} = \begin{pmatrix} |\lambda_1| & & \\ & |\lambda_2| & \\ & & \ddots \\ & & & |\lambda_n| \end{pmatrix}.$$

这里 D_1 是酉对角阵, $\tilde{\Lambda}$ 是对角元非负的对角阵, 类似地, 记

$$U_y = D_2 \tilde{U}_y, \quad I + R_x E_x R_x^{-1} = \tilde{Q}_x \tilde{R}_x. \quad (3.2.13)$$

这里 D_2 是酉对角阵, \tilde{U}_y 是对角元非负的上三角阵, \tilde{Q}_x, \tilde{R}_x 是 $I + R_x E_x R_x^{-1}$ 的 QR 分解, 于是

$$\begin{aligned} A' &= Q_x \tilde{Q}_x \tilde{R}_x R_x D_1' \tilde{\Lambda}' D_2 \tilde{U}_y \\ &= Q_x \tilde{Q}_x \tilde{R}_x R_x D_1' D_2 \tilde{\Lambda}' \tilde{U}_y \\ &= Q_x \tilde{Q}_x D_1' D_2 \tilde{R}_x \tilde{R}_x \tilde{\Lambda}' \tilde{U}_y. \end{aligned}$$

由于对角阵的乘法是可换的, 而酉对角阵右乘于对角元非负的上三角阵, 并从左边提出此酉对角阵时, 上三角阵中的对角元将保持不变, 即仍然为正, 故上式中 \tilde{R}_x, \tilde{R}_x 仍是对角元为正的上三角阵. 由 (3.2.13) 推知 \tilde{Q}_x, \tilde{R}_x 将收敛于单位阵, 并且 \tilde{R}_x 将收敛于单位阵. 由 A' 的 QR 分解的唯一性, 应有 $A' = P_1 U_1$ 及

$$P_1 = Q_x \tilde{Q}_x D_1' D_2, \quad U_1 = \Lambda_x \tilde{R}_x \tilde{\Lambda}' \tilde{U}_y, \quad (3.2.14)$$

又由

$$P_i = Q_1 \cdots Q_i = P_{i-1} Q_i,$$

有

$$\begin{aligned} Q_i &= P_{i-1}^{-1} P_i = D_2^{-1} D_1^{-i+1} \tilde{Q}_{i-1}^{-1} Q_{i-1}^{-1} Q_x \tilde{Q}_x D_1' D_2 \\ &= D_1^{-i+1} D_2^{-1} \tilde{Q}_{i-1}^{-1} \tilde{Q}_x D_1' D_2 \\ &= D_1^{-i+1} D_2^{-1} (I + F_{i-1}) D_1' D_2 \\ &= (I + \tilde{F}_{i-1}) D_1^{-i+1} D_2^{-1} D_1' D_2 \\ &= (I + \tilde{F}_{i-1}) D_1 \rightarrow D_1. \end{aligned} \quad (3.2.15)$$

这里 \tilde{F}_s 、 \tilde{R}_s 是根据 \tilde{Q}_s 收敛于单位阵而导出的收敛于零的阵。类似可证

$$\begin{aligned} R_s &= U_s U_s^{-1} \\ &= \tilde{R}_s \tilde{R}_s^{-1} \tilde{\Lambda}^s \tilde{U}_y \tilde{U}_y^{-1} \tilde{\Lambda}^{-s+1} \tilde{R}_s^{-1} \tilde{R}_s^{-1} \\ &= \tilde{R}_s \tilde{R}_s^{-1} \tilde{\Lambda} \tilde{R}_s^{-1} \tilde{R}_s^{-1} \rightarrow \tilde{R}_s \tilde{\Lambda} \tilde{R}_s^{-1}. \end{aligned}$$

故

$$A_s = Q_s R_s \rightarrow D_1 \tilde{R}_s \tilde{\Lambda} \tilde{R}_s^{-1} = AT \quad (3.2.16)$$

这里 T 是对角元为 1 的上三角阵。

2.2 等模特征值的情形

当条件 (3.2.7) 不满足时, 如果 A 只有线性初等因子, 这时 A 的 QR 分解序列的收敛性证明需要稍加修改, 得到的结论亦将有一些差别。设

$$|\lambda_1| > \cdots > |\lambda_r| = |\lambda_{r+1}| = \cdots = |\lambda_t| > |\lambda_{t+1}| > \cdots > |\lambda_n| > 0 \quad (3.2.17)$$

并且 Y 也有 LR 分解, 此时, 因

$$A^s = X \Lambda^s L_y R_y = X (\Lambda^s L_y \Lambda^{-s}) (\Lambda^s R_y),$$

其中 $\Lambda^s L_y \Lambda^{-s}$ 的位于区域

$$t \geq i > j \geq r \quad (3.2.18)$$

中的元素为

$$\begin{aligned} l_{ij} (\lambda_i / \lambda_j)^s &= l_{ij} \left(\frac{|\lambda_i| e^{i\theta_i}}{|\lambda_j| e^{i\theta_j}} \right)^s \\ &= l_{ij} e^{i(\theta_i - \theta_j)s}. \end{aligned} \quad (3.2.19)$$

它们的绝对值保持为 $|l_{ij}|$ ($t \geq i > j \geq r$), 而不随 $s \rightarrow \infty$ 而趋于零, 但 $\Lambda^s L_y \Lambda^{-s}$ 的在上述区域之外属于下三角阵的不在对角线上的元素则是趋于零的, 于是可写

$$A^s = X (\tilde{L}_s + E_s) (\Lambda^s R_y). \quad (3.2.20)$$

这里 \tilde{L}_s 仍是单位下三角阵, 其在 (3.2.18) 所示区域部分的元素相应地取为 $l_{ij}(\lambda_i/\lambda_j)^s$, 其余部分与单位阵相同, E_s 则是以 0 为极限的阵. 记

$$X\tilde{L}_s = \tilde{Q}_s\tilde{U}_s,$$

则

$$\begin{aligned} A^s &= \tilde{Q}_s\tilde{U}_s(I + \tilde{L}_s^{-1}E_s)\Lambda^s R_y \\ &= \tilde{Q}_s(I + \tilde{U}_s\tilde{L}_s^{-1}E_s\tilde{U}_s^{-1})\tilde{U}_s\Lambda^s R_y \\ &= \tilde{Q}_s(I + F_s)\tilde{U}_s\Lambda^s R_y \\ &= \tilde{Q}_s\bar{Q}_s\bar{R}_s\tilde{U}_s\Lambda^s R_y, \end{aligned} \quad (3.2.21)$$

其中 $F_s = \tilde{U}_s\tilde{L}_s^{-1}E_s\tilde{U}_s^{-1} \rightarrow 0$, $\bar{Q}_s\bar{R}_s$ 是 $I + F_s$ 的 QR 分解, 由 $F_s \rightarrow 0$, 可得 $\bar{Q}_s, \bar{R}_s \rightarrow I$, 注意, X 与 $X\tilde{L}_s$ 除第 r 列到第 t 列不同外, 其余都是相同的, $X\tilde{L}_s$ 的第 r 列到第 t 列分别是 X 的第 r 列到第 t 列的线性组合, 若令

$$X = Q_x R_x$$

为 X 的 QR 分解, 则由 $X\tilde{L}_s = \tilde{Q}_s\tilde{U}_s$ 及 $X\tilde{L}_s = Q_x R_x \tilde{L}_s$, 由于 $R_x \tilde{L}_s$ 具有形状

$$R_x \tilde{L}_s = \begin{pmatrix} \times & \times & \cdots & \times & \triangle & \cdots & \triangle & \times & \cdots & \times \\ & \times & & \vdots & \vdots & & \vdots & \vdots & & \vdots \\ & & & \times & \vdots & & \vdots & \vdots & & \vdots \\ & & & & \triangle & \cdots & \triangle & \vdots & & \vdots \\ & & & & \vdots & & \vdots & \vdots & & \vdots \\ & & & & \triangle & \cdots & \triangle & \times & & \vdots \\ & & & & & & & \times & & \vdots \\ & & & & & & & & \ddots & \vdots \end{pmatrix},$$

$\begin{matrix} r \\ t \end{matrix}$

其中画“ \times ”的元素与 R_x 中相应的元素相等, 画“ \triangle ”的元素则是 R_x 的第 t 列到第 r 列的元素的线性组合.

设 $R_x \tilde{L}_s$ 的 QR 分解为

$$R_x \tilde{L}_s = Q_A^{(s)} R_A^{(s)},$$

则由(3.2.22)的右端可知, $Q_{\mathcal{A}}^{(t)}$ 除第 t 到 r 行、列为一正交阵外, 其余与单位阵相同, 而 $R_{\mathcal{A}}^{(t)}$ 在 r 行上方各行和 t 行下方各行与 R, \tilde{L}_t 的相应的各行一致, 比较

$$X\tilde{L}_t = \tilde{Q}_t \tilde{U}_t = Q_x Q_{\mathcal{A}}^{(t)} R_{\mathcal{A}}^{(t)}, \quad (3.2.22)$$

可见

$$\tilde{Q}_t = Q_x Q_{\mathcal{A}}^{(t)}, \quad \tilde{U}_t = R_{\mathcal{A}}^{(t)}.$$

\tilde{Q}_t 与 Q_x 在第 r 列以左第 t 列以右是一样的, \tilde{Q}_t 的第 r 至第 t 列则是 Q_x 的第 r 至第 t 列的组合, 据(3.2.21)比较

$$A' = P_t U_t = \tilde{Q}_t \tilde{Q}_t^H \tilde{R}_t \tilde{U}_t A' R_t,$$

记 $A' = D' |A|'$, 其中 D' 是酉对角阵, 把上式 A' 中的酉对角阵 D' 往左移, 可见

$$P_t = \tilde{Q}_t \tilde{Q}_t^H D'.$$

据 $A_{s+1} = P_s^H A_s P_s = (D^H)^s \tilde{Q}_s^H \tilde{Q}_s^H A_s \tilde{Q}_s \tilde{Q}_s^H D_s$ 及(3.2.22)可得, $A_{s+1} = (D^H)^s \tilde{Q}_s^H \tilde{U}_s \tilde{L}_s^{-1} X^{-1} A_s X \tilde{L}_s \tilde{U}_s^{-1} \tilde{Q}_s D_s$, 由于 $\tilde{Q}_s \rightarrow I$, D_s 是酉对角阵, 故知 A_{s+1} 依实质收敛于

$$\tilde{U}_t (\tilde{L}_t^{-1} X^{-1} A_s X \tilde{L}_t) \tilde{U}_t^{-1}.$$

注意, $X\tilde{L}_t$ 的第 r 列以左和第 t 列以右应是 A_s 的对应于 $\lambda_1, \dots, \lambda_{r-1}, \lambda_{t+1}, \dots, \lambda_n$ 的特征向量, 于是 $\tilde{L}_t^{-1} X^{-1} A_s X \tilde{L}_t$ 有下列形状

$$\tilde{L}_t^{-1} X^{-1} A_s X \tilde{L}_t = \begin{pmatrix} \lambda_1 & & & & \\ & \ddots & & & \\ & & \lambda_{r-1} & & \\ & & & \times \cdots \times & \\ & & & \vdots & \\ & & & \times \cdots \times & \\ & & & & \lambda_{t+1} & \\ & & & & & \ddots \\ & & & & & & \lambda_n \end{pmatrix}, \quad (3.2.23)$$

\tilde{U}_s 与 \tilde{U}_s^{-1} 是互逆的上三角阵, 分别左乘和右乘于 (3.2.23) 后, 可知 A_{s+1} 依实质收敛于形如 (3.2.23) 的矩阵。

具有等模特征值的矩阵的一种特殊情况是, 矩阵 A 只具有单重实特征值和单重共轭复特征值。此时, A 的 QR 分解序列将依实质收敛于前述的拟上角阵。当矩阵有多个等模特征值分布在以原点为圆心的同心圆周上时, 特征值的确定将不如前面所述的情况直接, 这种困难将在后面予以克服。

定理 3.2.1 中要求 $X^{-1} \approx Y$ 有 LR 分解这个条件也可以去掉, 但在这种情况下, QR 分解的序列的极限矩阵的对角线上的元素可能不再按依模递减的次序来排列。有关的证明可参看 J. H. Wilkinson 的书的第 8 章 §31。

§3 带原点位移的 QR 算法

3.1 带位移的 QR 分解

当 A 的特征值依模不等时, 如果 Y 阵又有 LR 分解, 则 A 的 QR 分解序列将收敛于以 $\lambda_1, \lambda_2, \dots, \lambda_n$ 为对角元的上三角阵, A_s 的下三角部分的元素将以 $(\lambda_i/\lambda_j)^s (i > j)$ 的敛速收敛于零。由于 A 的特征值与 $A - pI$ 的特征值之差是常数 p , 而 $A - pI$ 的 QR 分解序列的下三角部分的元素将以 $(\lambda_i - p)^s / (\lambda_j - p)^s$ 的敛速收敛于零, 当 p 相对地更接近 λ_n 时, $(\lambda_n - p)^s / (\lambda_{n-1} - p)^s$ 收敛于零的速度就特别快。如果我们有一个确定 A 的依模最小的特征值 λ_n 的一个近似方法, 即若

$$k_s \rightarrow \lambda_n \quad (s \rightarrow \infty) \quad (3.3.1)$$

则 QR 分解过程必将得到加速。作

$$\begin{aligned} A_s - k_s I &= Q_s R_s, \quad A_{s+1} = R_s Q_s + k_s I \\ s &= 1, 2, \dots \end{aligned} \quad (3.3.2)$$

称 $\{A_s\}$ 为 A 的带恢复的以 k_s 为位移量的 QR 分解, 可以证明

$$A_{s+1} = R_s Q_s + k_s I = Q_s^H A_s Q_s, \\ A_{s+1} = Q_s^H \cdots Q_1^H A_1 Q_1 \cdots Q_s, \quad (3.3.3)$$

$$Q_1 \cdots Q_s A_{s+1} = A_1 Q_1 \cdots Q_s, \quad (3.3.4)$$

$$Q_1 \cdots Q_{s-1} (A_s - k_s I) R_{s-1} \cdots R_1 \\ = (A_1 - k_{s-1} I) (A_1 - k_{s-2} I) \cdots (A_1 - k_1 I). \quad (3.3.5)$$

又当适当选择 k_s 时, 则在 A_s 中, 其第 n 行将迅速出现收敛的情况, 当 $a_{ij}^{(s)}$ ($j < n$) 在计算精度要求下趋近于零时, 便可将 $a_{nn}^{(s)}$ 作为 A 的一个近似特征值输出, A 的其余特征值可在 A_s 的降维后的矩阵中寻找.

当 A 的依模最小的特征值是实的且是单重的时, 可取 $k_s = a_{nn}^{(s-1)}$. ($s = 1, 2, \cdots$).

当 A 的依模最小的特征值是共轭复数时, 可取 k_s 为 A_{s-1} 的右下角的二阶矩阵:

$$\begin{vmatrix} a_{n-1, n-1}^{(s-1)} & a_{n-1, n}^{(s-1)} \\ a_{n, n-1}^{(s-1)} & a_{n, n}^{(s-1)} \end{vmatrix}$$

的特征值, 这时 k_s 可能是复数. 为了在实际上使用复位移而又能避免复数运算, 需要进一步研究 QR 分解方法.

3.2 共轭复位移

设 k_1, k_2 是共轭的两个复数, 作

$$A_1 - k_1 I = Q_1 R_1, \quad A_2 = R_1 Q_1 + k_1 I, \quad (3.3.6)$$

$$A_2 - k_2 I = Q_2 R_2, \quad A_3 = R_2 Q_2 + k_2 I,$$

则 A_3 是 A_1 连续经 k_1, k_2 两次复位移后所得的矩阵, 当 A_1 是实阵时, 虽然 Q_1, R_1, Q_2, R_2 都是复阵, 但因

$$A_3 = (Q_1 Q_2)^H A_1 (Q_1 Q_2). \quad (3.3.7)$$

若令

$$Q = Q_1 Q_2, \quad R = R_2 R_1, \quad (3.3.8)$$

则有

$$\begin{aligned} QR &= Q_1 Q_2 R_2 R_1 = (A_1 - k_2 I)(A_1 - k_1 I) \\ &= A_1^2 - (k_1 + k_2)A_1 + k_1 k_2 I. \end{aligned} \quad (3.3.9)$$

上式的右端是实阵，它的 **QR** 分解也应是实阵，因而 Q 、 R 也是实的，由 (3.3.7)，从而 A_3 也是实的。

3.3 Hessenberg 阵

为了简化计算，我们证明下述定理。

定理 3.3.1 设 A 为实的非奇异阵， Q 是正交阵， B 是次对角元素为正数的实的上 Hessenberg 阵，则在等式

$$AQ = QB \quad (3.3.10)$$

中，当 Q 的第 1 列给定时， Q 、 B 将由 (3.3.10) 唯一确定。

证明 设 $Q = (q_1, \dots, q_n)$ 是正交阵，其中 q_1 是已知的，

$$B = \begin{pmatrix} b_{11} & b_{12} & \cdots & b_{1n} \\ & b_{22} & \cdots & b_{2n} \\ & & b_{33} & \cdots & b_{3n} \\ & & & \ddots & \vdots \\ & & & & b_{nn-1} & b_{nn} \end{pmatrix}, \quad (3.3.11)$$

比较 (3.3.10) 两端的第一列，得

$$Aq_1 = b_{11}q_1 + b_{21}q_2, \quad (3.3.12)$$

$$b_{11} = q_1^T Aq_1, \quad b_{21} = \|Aq_1 - b_{11}q_1\| > 0, \quad (3.3.13)$$

$$q_2 = (Aq_1 - b_{11}q_1) / b_{21}, \quad b_{21} \neq 0. \quad (3.3.14)$$

继续比较 (3.3.10) 两端的第 2, 3, ..., n 各列，并仿此讨论，知 B 、 Q 均被唯一确定：

$$\begin{aligned} b_{ik} &= q_i^T Aq_k \quad (i = k, \dots, n), \\ b_{k+1,k} &= \|Aq_k - \sum_{i=1}^k b_{ik}q_i\| > 0, \\ q_{k+1} &= (Aq_k - \sum_{i=1}^k b_{ik}q_i) / b_{k+1,k}, \\ k &= 1, 2, \dots, n-1. \end{aligned} \quad (3.3.15)$$

注意, 该定理把诸 $b_{k+1,k} \neq 0 (k = 2, \dots, n-1)$ 作为条件, 但此条件是否成立应在 (3.3.15) 的递推计算中才能发现, 一旦在某一步当 $b_{k+1,k} = 0$ 时, q_{k+1} 就不能再由 (3.3.15) 确定了, 此时, 可把 q_{k+1} 取作与 $q_i (i = 1, \dots, k)$ 正交的任意一个单位向量, 递推仍可继续下去, 但 Q 、 B 的唯一性就不再能保证了。

定理 3.3.2 设 A 是非奇异的上 Hessenberg 阵且有非零次对角元, 则 A 的 LR 和 QR 分解序列 $\{A_i\}$ 仍为上 Hessenberg 阵, 且有非零次对角元。

证明 设

$$A = LR \quad \text{或} \quad A = QU, \quad (3.3.16)$$

比较两端第 1 列, 有

$$a_1 = r_{11}l_1 \quad \text{或} \quad a_1 = u_{11}q_1. \quad (3.3.17)$$

由于 a_1 具有形状:

$$a_1^T = (x, x, 0, \dots, 0), \quad (3.3.18)$$

可见 l_1^T 、 q_1^T 亦具有 (3.3.18) 的形状, 依次比较 (3.3.16) 的其余各列, 可见 L 、 Q 亦均为上 Hessenberg 阵, 且有非零次对角元, 再由

$$A_2 = RL \quad \text{或} \quad \tilde{A}_2 = UQ$$

及 R 、 U 是上三角阵, 可知两个 A_2 仍是上 Hessenberg 阵, 同理 A_3, \dots, A_n, \dots 均是上 Hessenberg 阵。

由于任意矩阵均可经由 Householder 相似变换化为上 Hessenberg 阵, 当其次对角线上有零元时, 此种矩阵的特征值问题可以化为若干个低阶的次对角元非零的上 Hessenberg 阵的特征值问题。因此, 不失一般性, 今后我们将只讨论次对角元非零的上 Hessenberg 阵。它又被称为不可约的上 Hessenberg 阵。由于 Hessenberg 阵的 LR 或 QR 分解要比一般的满阵的分解在计算上省得多, 且又由于定理 3.3.2 的理由, 通常计算矩

阵的特征值时，总是先把给定矩阵化为上 Hessenberg 阵。在下面的讨论中，我们将始终假定 A_1 是不可约的上 Hessenberg 阵。

把定理 3.3.1 用到带两步共轭复位移的 QR 分解，由

$$\begin{cases} A_3 = Q^H A_1 Q, \quad Q = Q_1 Q_2, \\ A_1 Q = Q A_3, \end{cases} \quad (3.3.19)$$

可见，

$$A_1 \tilde{Q} = \quad (3.3.20)$$

中，当取 \tilde{Q} 的第 1 列为 Q 的第 1 列时，则由上式确定的 \tilde{Q} 、 B 便应分别是 Q 和 A_3 。

由 (3.3.9): $QR = A_1^3 - (k_1 + k_2)A_1 + k_1 k_2 I$ ，两端取第 1 列，得

$$r_{11} q_1 = (A_1 - k_1 I)(A_1 - k_2 I) e_1 \quad (3.3.21)$$

上式右端的列向量，在 A_1 是 Hessenberg 阵的假定下只有三个非零分量，即前三个分量，分别记为 x_1, y_1, z_1 则

$$\begin{aligned} x_1 &= a_{11}^3 + a_{12}a_{21} - a_{11}(k_1 + k_2) + k_1 k_2, \\ y_1 &= a_{21}(a_{11} + a_{22} - k_1 - k_2), \\ z_1 &= a_{21}a_{32}, \end{aligned} \quad (3.3.22)$$

于是

$$\begin{aligned} q_1^T &= \pm (x_1, y_1, z_1, 0, \dots, 0) / k, \\ k &= (x_1^2 + y_1^2 + z_1^2)^{1/2}. \end{aligned} \quad (3.3.23)$$

3.4 Francis 方法

有了 q_1 之后，本来就可按照 (3.3.12) ~ (3.3.15) 的公式计算 (3.3.19) 中的 Q 和 A_3 了。但是，按照这个过程去进行计算时，如果遇到上 Hessenberg 阵 A_3 的次对角元依模很小的情况，依 (3.3.15) 计算出的 $b_{k+1,k}$ 即 $a_{k+1,k}^{(3)}$ 和 q_{k+1} 的精度便很低，这将直接影响 Q 阵的正交性。因此，从计算稳定性看，这种算法是不可取的。Francis 对此做了改进，其计算过程可以归结为：求 $(n-1)$ 个 Householder 阵 P_1, \dots, P_{n-1} ，使 P_1 的第 1 列为 q_2 ，使 $P_1 \cdots P_k$ ， $(k = 1, 2, \dots, n-1)$ 的第 1 列也保持为 q_1 ，但

$$P_n^H \dots P_1^H A_1 P_1 \dots P_{n-1} \quad (3.3.24)$$

应是上 **Hessenberg** 阵。下面我们就来描述这些矩阵的构造方法。设

$$P_1 = I - 2\mathbf{w}_1\mathbf{w}_1^T, \quad (3.3.25)$$

由于 P_1 的第 1 列应是 q_1 ，而 q_1 只有前三个分量不为零，故 w_1 也是这样，令

$$\mathbf{w}_1 = \mathbf{v}_1 / \|\mathbf{v}_1\|, \quad \mathbf{v}_1^T = (1, \varphi_1, \varphi_2, 0, \dots, 0), \quad (3.3.26)$$

则

$$\|\mathbf{v}_1\|^2 = 1 + \varphi_1^2 + \varphi_2^2, \quad P_1 = I - 2\mathbf{v}_1\mathbf{v}_1^T / \|\mathbf{v}_1\|^2,$$

$$\begin{aligned} (P_1 \mathbf{e}_1)^T &= (1 - \frac{2}{\|\mathbf{v}_1\|^2}, -\frac{2\varphi_1}{\|\mathbf{v}_1\|^2}, -\frac{2\varphi_2}{\|\mathbf{v}_1\|^2}, 0, \dots, 0) \\ &= \pm (x_1, y_1, z_1, 0, \dots, 0)/k, \end{aligned} \quad (3.3.27)$$

故

$$\frac{2}{\|v_1\|^2} = -\frac{k+x_1}{k}. \quad (3.3.28)$$

为使上式中的分子不致相消, 令

$$\mathbf{k} = \text{diag}(x_1) (x_1^2 + y_1^2 + z_1^2)^{1/2}, \quad (3.3.29)$$

得

$$\varphi_1 = y_1 / (k + x_1), \quad \varphi_2 = z_1 / (k + x_1), \quad (3.3.30)$$

$$\|\mathbf{v}_1\|^2 = (k + x_1)/k; \quad (3.3.31)$$

于是 P_1 具有下列形状

$$P_1 = \left[\begin{array}{ccc|c} \times & \times & \times & \\ \times & \times & \times & \\ \times & \times & \times & 0 \\ \times & \times & \times & \\ \hline & & & \\ \hline & & & I_n \\ & & 0 & \end{array} \right]. \quad (3.3.32)$$

这是一个三阶 **Householder** 阵与一个单位阵的直接和。由于 $P_1 = P_1^H$ ，记 $B_1 = P_1 A P_1$ ，则 B_1 具有下列形状：

$$B_1 = \begin{pmatrix} \times & \times & \times & \cdots & \cdots & \cdots & \times \\ \times & \times & \times & & & & \times \\ \times & \times & \times & & & & \vdots \\ \times & \times & \times & \times & & & \vdots \\ & & & \times & \times & & \vdots \\ & & & & \ddots & & \vdots \\ & & & & & \ddots & \vdots \\ & & & & & & \times & \times \end{pmatrix}, \quad (3.3.33)$$

其第 4 行 4 列以后仍保持为上 **Hessenberg** 阵的形状。取 P_2 为：

$$P_2 = \begin{pmatrix} 1 & \vdots & 0 \\ \cdots & \cdots & \cdots \\ 0 & \vdots & U_2 \end{pmatrix}, \quad (3.3.34)$$

其中 U_2 为 $(n-1)$ 阶 **Householder** 阵，它把

$$B_1 = \begin{pmatrix} b_{11}^{(1)} & B_{12}^{(1)T} \\ \vdots & \vdots \\ b_{21}^{(1)} & \vdots \\ b_{31}^{(1)} & \vdots \\ b_{41}^{(1)} & \vdots \\ 0 & B_{22}^{(1)} \\ \vdots & \vdots \\ \vdots & \vdots \\ 0 & \vdots \end{pmatrix} = \begin{pmatrix} b_{11}^{(1)} & B_{12}^{(1)T} \\ \vdots & \vdots \\ B_{21}^{(1)} & B_{22}^{(1)} \\ \vdots & \vdots \\ \vdots & \vdots \end{pmatrix} \quad (3.3.35)$$

中的 $B_{21}^{(1)}$ 变形为

$$U_2 B_{21}^{(1)} = \begin{pmatrix} \bar{B}_{21}^{(1)} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad (3.3.36)$$

的形式。此时， U_2 实际上仍是一个三阶 **Householder** 阵与一个单位阵的直接和，而 $B_2 = P_2 B_1 P_2$ 则具有下列形状：

$$B_2 = P_2 B_1 P_2 = \begin{pmatrix} \times & \vdots & \times & \cdots & \times \\ \cdots & \vdots & \cdots & \cdots & \cdots \\ \times & \vdots & & & \\ 0 & \vdots & & & \\ \vdots & \vdots & & & \\ 0 & \vdots & & & \end{pmatrix} \quad (3.3.37)$$

$B_{22}^{(2)}$

它的第 1 列已经化为上 **Hessenberg** 型， $B_{22}^{(2)}$ 中除左上角在前三列含有一个 4×3 的矩阵外，其第 4 行 4 列以后仍保持为上 **Hessenberg** 型， $B_{22}^{(2)}$ 的形状与 B_1 的形状是一样的，因此，此后的计算便可组织在一个循环之中。每个 P_i 除 P_1 外，其任务是把 B_{i-1} 的第 i 列化成上 **Hessenberg** 型，而每个 P_i 的计算将归结为一个 U_i 的计算，其计算可仿 (3.3.25) ~ (3.3.31) 的计算进行。

这样，我们就有了一个计算矩阵 A 的带两步复共轭位移的 **QR** 分解的实的计算方法，其中分解阵 Q 具有较好的正交性。

另外，前曾指出，当矩阵有等模特征值时，**QR** 分解序列的极限将不再具有拟上三角阵的形式，极限阵中将出现阶数大于 2 的对角块，这种矩阵的特征值难于处理。在有了带位移特别是带两步共轭复位移的算法以后，矩阵有等模特征值的困难就不再严峻了。如图 3.1，设 A 为某实阵，其等模特征值分布在以原点为中心，以 r 为半径的某圆上且对称地出现在 x -轴的两侧。

当产生位移时，新原点如靠近某一特征值，必也靠近与其共轭的特征值，但同时却远离了其余的特征值，使原来有几对等模特征值的问题，在位移以后，化为只有一对最小的等模特征值问题。可见，使用带位移的QR算法，不仅能加速收敛过程，而且其适应能力也很强。

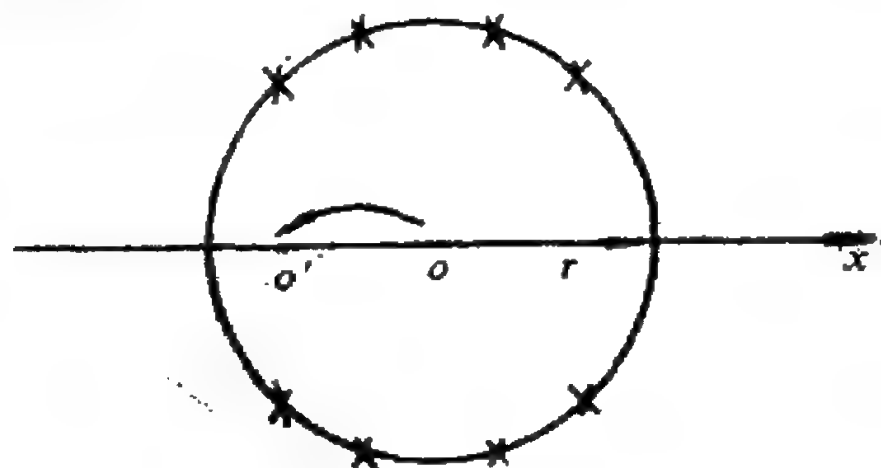


图 3.1 等模特征值

§4 QR 方法的细节

4.1 QR_2 子程序

用 QR_2 代表由 Hessenberg 矩阵 A 计算它的带两步位移量（实或共轭复数）的 QR 分解的实际计算程序，这个程序将是整个 QR 算法的核心子程序。作为其主要输入的形式参数包括：矩阵的阶数 n ，矩阵 A 和两个以实数形式给出的位移信息

$$\sigma = k_1 + k_2, \quad \rho = k_1 k_2, \quad (3.4.1)$$

其输出则是 (3.3.7) 中给出的 A_3 ，它将取代 A 而占有 A 原有的存贮单元。

考虑到 $\{A_i\}$ 始终保持着上 Hessenberg 阵的形式，其次对角元在工作精度内可能化为零。出现这种情况时，应分段地把

A_1 的特征值问题化为低阶阵的特征值问题，这样可以大大地节省计算工作量。为了实现这一点，可以从第 n 行往上数，把合于

$$|a_{i,i-1}| < \varepsilon \quad (i = n, n-1, \dots, 1) \quad (3.4.2)$$

的最大足码记录下来，设为 n_p ，QR分解序列计算便只在第 n_p 行（列）到第 n 行（列）之间进行。

在实现

$$A_3 = B_{n-1} = P_{n-1} \cdots P_1 A_1 P_1 \cdots P_{n-1}$$

的计算时，可按递推公式

$$B_j = P_j B_{j-1} P_j, \quad j = 1, \dots, n-1$$

$$B_0 = A_1$$

去组织循环，由

$$P_j = I - v_j (\alpha v_j^T),$$

$$\alpha = \left(1 + \frac{x_1}{k}\right),$$

$$k = \text{sign}(x_1) (x_1^2 + y_1^2 + z_1^2)^{1/2}, \quad (3.4.3)$$

$$v_j = e_j + \varphi_1 e_{j+1} + \varphi_2 e_{j+2},$$

$$\varphi_1 = y_1 / (x_1 + k), \quad \varphi_2 = z_1 / (x_1 + k);$$

则当 x_1, y_1, z_1 给定时，由 A_1 计算 A_3 的过程就完全可以进行了。令

$$B_j = B'_{j-1} P_j, \quad (3.4.4)$$

$$B'_{j-1} = P_j B_{j-1} = (I - v_j (\alpha v_j^T)) B_{j-1}, \quad (3.4.5)$$

$$\eta = \alpha v_j^T B_{j-1}; \quad (3.4.6)$$

则

$$B'_{j-1} = B_{j-1} - v_j \eta^T. \quad (3.4.7)$$

于是

$$\begin{aligned} B_j &= B'_{j-1} - \xi v_j^T, \\ \xi &= \alpha B'_{j-1} v_j. \end{aligned} \quad (3.4.8)$$

计算每个 P_i 所需的输入信息 x_i, y_i, z_i 只有一点不同, 即 P_1 所需的 x_1, y_1, z_1 是由

$$x_1 = a_{12}^2 + a_{12}a_{21} - \sigma a_{11} + \rho,$$

$$y_1 = a_{21}a_{11} + a_{22}a_{21} - \sigma a_{21},$$

$$z_1 = a_{32}a_{21}$$

计算出来的, (参看 (3.3.22)), 而 P_2, \dots, P_{n-1} 所需的 x_i, y_i, z_i 则应分别由 B_1, \dots, B_{n-2} 中的元素提供, 即

$$x_i = b_{j+1,j}^{(j)},$$

$$y_i = b_{j+2,j}^{(j)}, \quad j = 1, 2, \dots, n-2, \quad (3.4.9)$$

$$z_i = b_{j+3,j}^{(j)}, \quad j+3 > n \text{ 时 } z_i = 0.$$

于是有关 B_1, B_2, \dots, B_{n-1} 的计算几乎可以在同一循环下统一地进行处理了。这样编程的结果, 从 A_1 算出 A_3 大约只需 $5n^2$ 次乘法, 而在降阶的情况下还会更省一些。

4.2 位移的选择

使用带位移的 QR 分解技术, 目的在于加速 QR 分解序列的收敛进程。然而, 在当过程显露出某种收敛趋势时, 这种技术才是比较有把握的。往往有一些简单的措施, 虽不曾在理论上证明是最好的, 但却能帮助我们去“捕捉”这种趋势。

把矩阵的右下角的二阶矩阵的特征值 μ_1 和 μ_2 计算并存放起来, 称之为老的近似特征值。当 μ_1, μ_2 是复数时, 应分别把它们的实部和虚部存放起来。然后, 用不带位移的两步 QR 分解法求 A_3 并也求出新的 μ_1 和 μ_2 。如果把老的 μ_1, μ_2 和新的 μ_1, μ_2 分别记为 $\mu_1^{(j-1)}, \mu_2^{(j-1)}$ 和 $\mu_1^{(j)}, \mu_2^{(j)}$, 并约定当两个 μ 均为实数时, 把数值较小的那个放到 μ_1 , 另一个放到 μ_2 , 当两个 μ 是共轭复数时, 把有负虚部的那个当作 μ_1 而把有正虚部的当作 μ_2 。当

$$|\mu_i^{(r)} - \mu_i^{(r-1)}| > \frac{1}{2} |\mu_i^{(r)}|, \quad (i=1, 2) \quad (3.4.10)$$

时，即认为过程还未显露收敛的趋势。这时可令

$$k_i^{(r)} = k_i^{(r-1)} \quad (i=1, 2). \quad (3.4.11)$$

这表示强令新老位移一致，即仍用原来的位移，而 $k_i^{(0)} = 0$, $(i=1, 2)$ ，即在开始时的位移量是取为零的，这里用 k_1, k_2 作存放位移的单元，而 $k_i^{(r)} (i=1, 2)$ 是其流动值。当

$$|\mu_i^{(r)} - \mu_i^{(r-1)}| \leq \frac{1}{2} |\mu_i^{(r)}|, \quad (i=1, 2) \quad (3.4.12)$$

时可取

$$k_i^{(r)} = \mu_i^{(r)} \quad (i=1, 2), \quad (3.4.13)$$

即用新算出的近似特征值作为位移量，当 (3.4.10) 只有一个成立时，由于使 (3.4.12) 成立的那个 $\mu_i^{(r)} (i=1 \text{ 或 } 2)$ ，可能是实数也可能是复数，当其是实数时，由于另一个 μ 应当相弃，故应取

$$k_1^{(r)} = k_2^{(r)} = \mu_i^{(r)} \quad (i=1 \text{ 或 } 2),$$

从而

$$\sigma = 2\mu_i^{(r)}, \quad \rho = (\mu_i^{(r)})^2.$$

当 $\mu_i^{(r)}$ 是复数时，可取

$$k_1^{(r)} = k_2^{(r)} = \operatorname{Re}(\mu_i^{(r)}),$$

得

$$\sigma = 2\operatorname{Re}(\mu_i^{(r)}), \quad \rho = (\operatorname{Re}(\mu_i^{(r)}))^2. \quad (3.4.14)$$

(3.4.14) 对 $\mu_i^{(r)} (i=1 \text{ 或 } 2)$ 是实数或复数都是通用的。因此，当 (3.4.12) 中只有一个成立时，我们就用 (3.4.14) 去确定新位移的信息。(3.4.10) 中的 $\frac{1}{2}$ 因子取得稍大或稍小一点是无关紧要的。

当矩阵有等模特征值时，如等模特征值个数 r 大于 2，条件 (3.4.12) 可能经多次试验都不满足，正如 QR 方法收敛性证明中所作的分析那样， A_s 中将有一个 r 阶的对角块，其元素是摆动而不收敛的，当所论特征值的模最小时，这个对角块将出现在 A_s 的右下角。为了克服这种情况，可在程序中设置一个记录使 (3.4.12) 不成立的计数器，当使 (3.4.12) 不成立的次数超过某一规定的正整数 T_1 时，就把矩阵归结为有等模特征值且又是个数大于 2 的情况。由于矩阵的行列式的绝对值等于矩阵的各特征值的模的乘积。对于等模情况，这个模应等于 $|\det(C)|^{1/r}$ 。因此，位移量可取为

$$k_1 = k_2 = \pm |\det(C)|^{1/r}. \quad (3.4.15)$$

这里 C 是 A_s 的由第 n_p 行（列）到第 n 行（列）组成的矩阵， r 是它的阶，上式中“ \pm ”号的选择可使之与 $\text{trace}(C)$ 同号。 T_1 的大小建议取为

$$T_1 = \min(n - n_p + 1, 10).$$

4.3 工作精度与矩阵的预处理

通常可以用矩阵 A 的欧氏范数 $\|A\|_E$ 作为 A 的特征值的上界，如果机器字长为 t ，则此时合理的工作精度应取为

$$\varepsilon = 2^{-t} \|A\|_E, \quad \|A\|_E^2 = \sum_{i,j} a_{ij}^2. \quad (3.4.16)$$

但在一些实际问题中， $\|A\|_E$ 可能比 A 的依模最大的特征值还大许多，这时，按 (3.4.16) 来确定 ε 就不合实际了。这时，应利用对角相似变换把 $\|A\|_E$ 适当减小，这个工作叫做对矩阵进行平衡化，其要求是使矩阵的各行各列中依模最大的元素，不要在数量级上有太大的差别。

例 3.4.1 设

$$A = \begin{pmatrix} 1 & 100 & & \\ 1 & 1 & 1 & \\ & 100 & 1 & 100 \\ & & & 1 \end{pmatrix},$$

于是 $\|A\|_E = \sqrt{30007} = 173.2$, 如令 $D = \text{diag}(10, 1, 10, 1)$ 则有

$$\tilde{A} = D^{-1}AD = \begin{pmatrix} 1 & 10 & & \\ 10 & 1 & 10 & \\ & 10 & 1 & 10 \\ & & & 10 & 1 \end{pmatrix}.$$

从而 $\|\tilde{A}\|_E = \sqrt{604} = 24.5$, 矩阵的实际依模最大的特征值为 17.1.

4.4 计算次序

用 QR 方法求实矩阵的全部工作内容(包括预处理)如下:

- (i) 利用对角相似变换使 $\|A\|_E$ 适当减小;
- (ii) 计算 $\|A\|_E$, $\text{trace}(A)$ 及 $\varepsilon = 2^{-i} \|A\|_E$;
- (iii) 用正交相似变换使 A 化为上 Hessenberg 阵.
- (iv) 计算右下角的二阶矩阵的特征值 μ_1 及 μ_2 , 并给定初始位移量 $k_1 = k_2 = 0$;
- (v) 求 Hessenberg 阵的次对角元为零元的位置, 确定 n_p , 当 $n_p \geq n-1$ 时输出 A 的一个或两个特征值, 并使矩阵的阶数减 1 或 2;
- (vi) 调用一次 QR_2 程序;
- (vii) 重新计算 μ_1 及 μ_2 , 并依 (3.4.10) ~ (3.4.14) 重新修改位移量, 必要时依 (3.4.15) 确定位移量;
- (viii) 重复 (v) ~ (vii) 的循环, 直到算出全部的特征值为止.

例 3.4.2 设

$$A_1 = \begin{pmatrix} 5.0 & -2.0 & -5.0 & -1.0 \\ 1.0 & 0.0 & -3.0 & 2.0 \\ & 2.0 & 2.0 & -3.0 \\ & & 1.0 & -2.0 \end{pmatrix}$$

算得 $\mu_1^{(1)} = -1$, $\mu_2^{(1)} = 1$, 开始时置 $\sigma_1 = \rho_1 = 0$, 调用一次 QR₂, 得

$$A_3 = \begin{pmatrix} 4.11828 & 3.46451 & -5.39263 & 1.51084 \\ -0.29992 & 0.68721 & -3.30064 & -4.23760 \\ & 0.93016 & 1.11718 & -0.71200 \\ & & 0.22015 & -0.92266 \end{pmatrix}$$

由 A_3 算得 $\mu_1^{(3)} = -0.84268$, $\mu_2^{(3)} = 1.03720$, 此时条件 (3.4.12) 均满足, 应取 $\sigma_3 = \mu_1^{(3)} + \mu_2^{(3)} = 0.19452$, $\rho_3 = \mu_1^{(3)} \mu_2^{(3)} = -0.87403$. 经调用 QR₂, 得

$$A_5 = \begin{pmatrix} 3.88102 & 5.50657 & -2.83249 & -9.8166 \\ 0.08301 & 2.29698 & -3.41894 & 3.17136 \\ & 1.65698 & -0.18116 & 2.41387 \\ & & 0.01060 & -0.99684 \end{pmatrix}$$

由 A_5 得 $\mu_1^{(5)} = -1.02709$, $\mu_2^{(5)} = -0.15091$. 此时条件 (3.4.12) 只对 $\mu_1^{(5)}$ 成立, 应取 $\sigma_5 = 2\mu_1^{(5)} = -2.05418$, $\rho_5 = (\mu_1^{(5)})^2 = 1.05491$,

再调 QR₂ 得

$$A_7 = \begin{pmatrix} 4.01900 & 0.73616 & -6.17208 & 0.85558 \\ -0.02762 & -0.19762 & -3.26500 & -3.83089 \\ & 1.68124 & 2.17861 & 1.23705 \\ & & -\varepsilon & -1.00000 \end{pmatrix}$$

由 A_7 得 $\mu_1^{(7)} = -1$, $\mu_2^{(7)} = 2.17861$. 此时条件 (3.4.12) 也只对 $\mu_1^{(7)}$ 成立, 应取 $\sigma_1 = 2\mu_1^{(7)} = -2$, $\rho_1 = (\mu_1^{(7)})^2 = 1$, 得

$$A_9 = \begin{pmatrix} 4.01231 & 5.45684 & -3.02039 & -0.83340 \\ -0.00903 & 2.11330 & -3.33301 & 3.23455 \\ & 1.57244 & 0.12561 & 2.40441 \\ & & -0. & -1.00000 \end{pmatrix}.$$

此时 (4,3) 元素已经在机器精度内化为零, 一个特征值 $\lambda_4 = -1$, 已经可以输出, 去掉第 4 行 4 列后, 仍取初始位移信息 $\sigma_7 = \rho_7 = 0$, 计算在第 1 至 3 行 3 列进行, 算得

$$\mu_1^{(9)} = 1.119455 - 2.063361i,$$

$$\mu_2^{(9)} = 1.119455 + 2.063361i,$$

调用 QR_2 后得

$$A_{11} = \begin{pmatrix} 3.99618 & 4.15400 & 4.64849 & \vdots & * \\ 0.99182 & 0.37161 & 1.17200 & \vdots & * \\ & -3.74868 & 1.63220 & \vdots & * \\ \dots\dots\dots & & & \vdots & \dots\dots\dots \\ & & & \vdots & -1 \end{pmatrix}.$$

由 A_{11} 得 $\mu_1^{(11)} = 1.001905 - 1.99905i$, $\mu_2^{(11)} = 1.001905 + 1.99905i$. 此时条件 (3.4.12) 对两个 μ 都满足, 应取 $\sigma_{11} = \mu_1^{(11)} + \mu_2^{(11)} = 2.00381$, $\rho_{11} = \mu_1^{(11)}\mu_2^{(11)} = 5.0000$. 再算得

$$A_{13} = \begin{pmatrix} 4.00000 & 5.04835 & -3.65643 & \vdots & * \\ 0 & 1.87849 & -3.59100 & \vdots & * \\ & 1.32902 & 0.12106 & \vdots & * \\ \dots\dots\dots & & & \vdots & \dots\dots\dots \\ & & & \vdots & -1.0 \end{pmatrix}.$$

由 A_{13} 得 $\mu_1^{(13)} = 1.0 - 2.0i$, $\mu_2^{(13)} = 1.0 + 2.0i$, 于是求得全部特征值为: $\lambda_1 = 4$, $\lambda_{2,3} = 1.0 \mp 2.0i$, $\lambda_4 = -1$.

§ 5 实对称阵的 QR 算法

5.1 QR 算法与对称性

把 QR 方法应用于对称矩阵, 理应得出更进一步的结果。我们首先指出, 对称性在 QR 变换下仍是保持的。设 $A = A_1$ 是实对称阵, A_2 是 A_1 经一次 QR 分解后得到的矩阵:

$$A_1 = QR, \quad A_2 = RQ, \quad (3.5.1)$$

由于 $A_2 = Q^T A_1 Q$ 及 A_1 的对称性可知, A_2 也是对称的。

根据这个性质, 当把 QR 方法应用于对称阵时, 矩阵 A 及 A_1 均可只存贮其上三角阵部分, 这样存贮量与计算量可大致节约一半。

然而, LR 方法却不具有这种性质。但若把 LR 分解修改成 Cholesky 分解, 却仍能保持对称性。即令

$$A_1 = LL^T, \quad A_2 = L^T L, \quad (3.5.2)$$

则

$$A_2 = L^T A_1 L^{-T}. \quad (3.5.3)$$

可见, A_2 仍是对称的。但在此时, 不仅要求 A_1 是对称的, 而且要求 A_1 是正定的, 否则 (3.5.2) 的分解中将遇到复数运算。这是 QR 方法应用于对称阵时所不具有而为 LR 方法所具有的缺点。

把 QR 方法与 Cholesky 方法加以比较, 可进一步看出二者的差别和联系。由于

$$A_1^k = (Q_1 \cdots Q_k)(R_1 \cdots R_k) \quad (3.5.4)$$

及

$$\begin{aligned} A_1^{k+1} &= (A_1^k)^T A_1^k = (R_1 \cdots R_k)^T (Q_1 \cdots Q_k)^T (Q_1 \cdots Q_k) (R_1 \cdots R_k) \\ &= (R_1 \cdots R_k)^T (R_1 \cdots R_k), \end{aligned} \quad (3.5.5)$$

而由 A_1 的 Cholesky 分解, 有

$$\begin{cases} \tilde{A}_s = L_s L_s^T, & \tilde{A}_{s+1} = L_s^T L_s, \quad s = 1, 2, \dots, \\ \tilde{A}_1 = A_1 \end{cases} \quad (3.5.6)$$

可得

$$A_1^{2s} = (L_1 \cdots L_{2s}) (L_1 \cdots L_{2s})^T \quad (3.5.7)$$

和

$$A_{s+1} = \tilde{A}_{2s+1}. \quad (3.5.8)$$

可见, A_1 的 Cholesky 分解序列的第 $2s+1$ 个就是 A_1 的 QR 分解序列的第 $s+1$ 个. 由于实对称阵的 QR 分解必是实的, 于是 (3.5.8) 又是任意实对称阵的奇数次 Cholesky 分解必是实的一个证明.

由于 QR 分解的序列将不改变矩阵的 2-范数和欧氏范数, 矩阵序列中各元素显然是有界的, 而在 Cholesky 分解中, 当矩阵近于奇异时, 其三角阵因子中的元素的上界, 一般是无法估计的, 所以从算法稳定性来看, 这又是 QR 方法的另一优点.

5.2 QR 算法的三阶收敛性

设 $A = A_1$ 是对称阵, $\{A_s\}$ 是由 A_1 产生的 QR 分解序列, 记

$$A_s = \begin{bmatrix} A_{11}^{(s)} & A_{12}^{(s)} \\ A_{21}^{(s)} & A_{22}^{(s)} \end{bmatrix} \quad (3.5.9)$$

算法的收敛速度将用 $\|A_{11}^{(s)}\|_E \rightarrow 0$ 的速度来衡量, 即如能由 $\|A_{11}^{(s)}\|_E = |\varepsilon|$, 导出

$$\|A_{21}^{(s+1)}\|_E = O(|\varepsilon|^p) \quad (p > 1) \quad (3.5.10)$$

对充分大的 s 成立, 则称算法是 p -阶收敛的. 这里 $A_{11}^{(s)}$ 是 r 阶方阵, 而 $1 \leq r \leq n$.

定理 3.5.1 由实对称阵 A_1 生成的带位移的 QR 分解序列具有三阶的敛速, 而位移量的选择始终用 $k_s = a_{nn}^{(s)}$.

证明 设 A_1 的依模最小的特征值 λ_n 的重数是 r , 即

$$|\lambda_1| \geq \cdots \geq |\lambda_{n-r}| > |\lambda_{n-r+1}| = \cdots = |\lambda_n|, \quad (3.5.11)$$

并且 $\lambda_{n-r+1} = \cdots = \lambda_n$, 记

$$A_s = \begin{bmatrix} F_s & G_s \\ G_s^T & H_s \end{bmatrix}^{(n-r)}. \quad (3.5.12)$$

据 QR 方法的收敛性的一般理论可知, G_s 将以零阵为极限, 记 $\lambda'_i (i=1, 2, \dots, n-r)$ 为 F_s 的特征值, $\lambda''_i (i=n-r+1, \dots, n)$ 为 H_s 的特征值, 当 s 充分大时, 可设

$$A_s = \begin{bmatrix} F_s & \varepsilon k_s \\ \varepsilon k_s^T & H_s \end{bmatrix}, \quad \|k_s\|_E = 1 \quad (3.5.13)$$

及

$$|\lambda'_i - \lambda_n| > \frac{2}{3} (|\lambda_{n-r} - \lambda_n|), \quad (i=1, \dots, n-r), \quad (3.5.14)$$

$$\varepsilon < \frac{1}{3} |\lambda_{n-r} - \lambda_n|, \quad (3.5.15)$$

由于 A_s 以 λ_n 为 r 重特征值, 故 $A_s - \lambda_n I$ 的秩为 $(n-r)$. 另一方面, $F_s - \lambda_n I$ 的 $n-r$ 个特征值为 $\lambda'_i - \lambda_n$ ($i=1, \dots, n-r$), 它们均不为零, 故 $F_s - \lambda_n I$ 的秩也是 $(n-r)$, 设 $F_s - \lambda_n I$ 有三角分解

$$F_s - \lambda_n I = L_s R_s, \quad (3.5.16)$$

作

$$L = \begin{bmatrix} L_s^{-1} & 0 \\ P_s & I \end{bmatrix}^{(n-r)}, \quad (3.5.17)$$

其中 P_s 满足

$$P_s (F_s - \lambda_n I) + \varepsilon k_s^T = 0. \quad (3.5.18)$$

由于 $F_s - \lambda_n I$ 有逆, P_s 显然有解, 于是有

$$L(A_s - \lambda_n I) = \begin{bmatrix} R_s & \varepsilon L_s^{-1} k_s \\ 0 & \varepsilon P_s k_s + H_s - \lambda_n I \end{bmatrix}, \quad (3.5.19)$$

上式两端的秩均应是 $n-r$, 必有

$$\varepsilon P_s k_s + H_s - \lambda_n I = 0 \quad (3.5.20)$$

与 (3.5.18) 联立消去 P_s 得

$$-\varepsilon^2 k_s^T (F_s - \lambda_n I)^{-1} k_s + (H_s - \lambda_n I) = 0, \quad (3.5.21)$$

$$\begin{aligned} H_s &= \lambda_n I + \varepsilon^2 k_s^T (F_s - \lambda_n I)^{-1} k_s \\ &= \lambda_n I + M_s, \end{aligned} \quad (3.5.22)$$

而

$$M_s = \varepsilon^2 k_s^T (F_s - \lambda_n I)^{-1} k_s, \quad (3.5.23)$$

于是

$$\begin{aligned} \|M_s\|_E &= \varepsilon^2 \|k_s^T (F_s - \lambda_n I)^{-1} k_s\|_E \\ &\leq \varepsilon^2 \|k_s^T (F_s - \lambda_n I)^{-1}\|_E \|k_s\|_E \\ &\leq \varepsilon^2 \|k_s^T\|_E \|(F_s - \lambda_n I)^{-1}\|_2^* \\ &\leq \varepsilon^2 / \left[\frac{2}{3} |\lambda_{n-r} - \lambda_n| \right] \\ &\leq \frac{3}{2} \cdot \frac{1}{9} |\lambda_{n-r} - \lambda_n| = \frac{1}{6} |\lambda_{n-r} - \lambda_n|. \end{aligned} \quad (3.5.24)$$

特别, 有

$$|a_{nn}^{(s)} - \lambda_n| < \frac{3}{2} \varepsilon^2 / |\lambda_{n-r} - \lambda_n| < \frac{1}{6} |\lambda_{n-r} - \lambda_n|. \quad (3.5.25)$$

据此, 可导出

$$|\lambda_i' - a_{nn}^{(s)}| > |\lambda_i' - \lambda_n| - |\lambda_n - a_{nn}^{(s)}| > \frac{1}{2} |\lambda_{n-r} - \lambda_n| \quad (3.5.26)$$

$$(i = 1, 2, \dots, n-r),$$

$$|\lambda_i'' - a_{nn}^{(s)}| \leq |\lambda_i'' - \lambda_n| + |\lambda_n - a_{nn}^{(s)}| < 3\varepsilon^2 / |\lambda_{n-r} - \lambda_n| \quad (3.5.27)$$

$$(i = n-r+1, \dots, n).$$

下面我们研究对 A_s 实行一步带位移的QR算法的效果. 位移量取为 $a_{nn}^{(s)}$, 令

$$P_s = \begin{bmatrix} Q_s & R_s \\ S_s & T_s \end{bmatrix}, \quad (3.5.28)$$

以使 $P_s(A_s - a_{nn}^{(s)}I)$ 成为上三角阵的正交阵, 于是

* 此处用到不等式: $\|AB\|_2 \leq \min\{\|A\|_2 \|B\|_E, \|A\|_E \|B\|_2\}$.

$$\begin{aligned} A_{i+1} &= P_i (A_i - a_{nn}^{(s)} I) P_i^T \\ &= \begin{bmatrix} F_{i+1} & G_{i+1} \\ G_{i+1}^T & H_{i+1} \end{bmatrix}, \end{aligned}$$

应有

$$S_i (F_i - a_{nn}^{(s)} I) + \varepsilon T_i k_i^T = 0, \quad (3.5.29)$$

$$G_{i+1}^T = (\varepsilon S_i k_i + T_i (H_i - a_{nn}^{(s)} I)) R_i^T. \quad (3.5.30)$$

利用 $\|R_i\|_E = \|S_i\|_E$, 可得

$$\|G_{i+1}^T\|_E \leq \varepsilon \|S_i\|_E^2 + \|T_i\|_E \|S_i\|_E \|H_i - a_{nn}^{(s)} I\|_2$$

及

$$\begin{aligned} \|S_i\| &= \varepsilon \|T_i k_i^T (F_i - a_{nn}^{(s)} I)^{-1}\|_E \\ &\leq \varepsilon \|T_i\|_E \max_{1 \leq i \leq n-r} |(\lambda_i^{-1} - a_{nn}^{(s)})^{-1}| \\ &< 2\varepsilon \|T_i\|_E / |\lambda_{n-r} - \lambda_n|. \end{aligned} \quad (3.3.31)$$

最后

$$\begin{aligned} \|G_{i+1}\|_E &\leq 4\varepsilon^3 \|T_i\|_E^2 / |\lambda_{n-r} - \lambda_n|^2 \\ &\quad + 2\varepsilon \|T_i\|_E^2 \max |\lambda_i'' - a_{nn}^{(s)}| / |\lambda_{n-r} - \lambda_n| \\ &< 10\varepsilon^3 \|T_i\|_E^2 / |\lambda_{n-r} - \lambda_n|^2 \\ &\leq 10r\varepsilon^3 / |\lambda_{n-r} - \lambda_n|^2. \end{aligned}$$

这就证明了算法的三阶收敛性。由证明过程还可看出, A 的依模最小的 r 个特征值是同时求出的, 于是有重特征值的情形对于对称矩阵来说, 就不再是一种缺点而是一种优点了。

5.3 对称带状阵

设 A 是带宽为 $2m+1$ 的实对称阵, 即

$$\begin{aligned} a_{ij} &= a_{ji} & i, j &= 1, \dots, n, \\ a_{ij} &= 0 & |i-j| &> m, \end{aligned} \quad (3.5.32)$$

首先指出: 对称带状阵在 QR 分解变换后仍为对称带状阵, 且带宽不变。

令 $P_i (i = 1, 2, \dots, n-1)$ 是初等 Hermite 阵, 其中 P_i 将使 $P_i(P_{i-1} \cdots P_1 A)$ 的第 i 列上三角化, 且 $P_i \cdots P_1 A$ 与 $P_{i-1} \cdots P_1 A$ 的前 $i-1$ 列是一样的, 于是

$$P_{n-1} \cdots P_1 A = R \quad (3.5.33)$$

是上三角阵, 它的上半带宽是 $2m+1$ 而不是 $m+1$, 令

$$Q = P_1 \cdots P_{n-1},$$

有

$$A = QR$$

及 $A_2 = RQ = RP_1 \cdots P_{n-1}$, 依次将 P_1, P_2, \dots, P_{n-1} 自右边连乘于 R , 可见 $RP_1 \cdots P_{n-1}$ 的下三角部分是带宽为 $m+1$ 的阵, 由于 $RP_1 \cdots P_{n-1} = QR$ 应是对称阵, 故 QR 实际是带宽为 $2m+1$ 的对称阵。

在讨论对称带状阵的带位移的 QR 分解时, 位移量 k_i 的选择, 仍可按 §4 中所述的一般原则进行。本段所要讨论的是节约内存的最大可能性。

由于 A_1 是对称带状阵, 可以只输入其上三角部分。为输入方便, 可以把 A 的非零元 $a_{i, k+i}$ 放在二维数组 $B[1:n; 0:m]$ 的 $B[i, k]$ 单元中, 这里 $i = 1, 2, \dots, n, k = 0, 1, \dots, \min(n-i, m)$ 。由 A_1 计算 A_2 时, 主要经过以下两个步骤:

(i) 由 A_1 计算 $R = P_{n-1} \cdots P_1 A_1$;

(ii) 由 R 计算 $A_2 = RP_1 \cdots P_{n-1}$ 。

由于 $P_i P_{i-1} \cdots P_1 A_1$ 不改变 $P_{i-1} \cdots P_1 A$ 的前 $i-1$ 行, 所以 $P_i \cdots P_1 A_1$ 的第 i 行的输出就是 $P = P_{n-1} \cdots P_1 A_1$ 的第 i 行的输出, 由于 $P_i \cdots P_1 A$ 的第 i 行的宽度是 $2m+1$ 而不是 $m+1$, 这样, 用它去冲掉 A_1 的第 i 行是有困难的。但是如果注意到在由 R 计算 A_2 时, 有些不必记录的中间结果可以及时剔除掉, 以便尽可能的节省存贮量和计算量。

以 $m=2, n=7$ 为例, 我们将通过比较 R 和 RP_1 的存贮情况来描述这个问题的计算与存贮的策略

R	RP_1
$\begin{pmatrix} \times & \times & \times & \times & \times & \bigcirc & \bigcirc \\ & \times & \times & \times & \times & \times & \bigcirc \\ & & \times & \times & \times & \times & \times \\ & & & \times & \times & \times & \times \\ & & & & \times & \times & \times \\ & & & & & \times & \times \\ & & & & & & \times \end{pmatrix}$	$\begin{pmatrix} \triangle & \bullet & \bullet & \times & \times & & \\ \triangle & \otimes & \otimes & \times & \times & \times & \\ \triangle & \otimes & \otimes & \times & \times & \times & \times \\ & & & \times & \times & \times & \times \\ & & & & \times & \times & \times \\ & & & & & \times & \times \\ & & & & & & \times \end{pmatrix}$

R 在右乘以 P_1 后, 它只影响 R 的前三列, 即画有 \triangle 、 \bullet 和 \otimes 的那些元素。注意, 画有 \triangle 的元素, 在以后继续右乘 P_2, \dots 时将不再变化, 即它们是 A_2 的第 1 列的最后输出信息; 画有 \otimes 的元素在 RP_1 继续右乘 P_2, \dots 时, 它们将发生变化并且对于形成 A_2 的第 2、3 列在对角线下的元素是有贡献的; 画有 \bullet 的元素, 在以后的计算中, 虽然也要发生变化, 但它们对形成 A_2 的对角线以下的第 2 列以后的元素却没有贡献。仿此分析可知, 在 R 中真正对形成 A_2 的各列在对角线以下的元素有贡献的只是 R 中位于第 1, 2, $\dots, m+1$ 条对角线上的那些元素。可见 R 在右乘以 P_1 之后, RP_1 的第 1 行中的前 $m+1$ 个元素可以用 RP_1 的第 1 列中画有 \triangle 的元素冲掉。同样, 当右乘 P_2 以后, 可用 RP_1P_2 的第 2 列在对角线以下的元素冲掉 RP_1P_2 的第 2 行对角线以右的那些元素。注意, 当

$$P_{m+1} \cdots P_1 A_1$$

算出以后, 前 $m+1$ 列便已经上三角化, 而后继的流动矩阵

$$P_k \cdots P_{m+2} (P_{m+1} \cdots P_1 A_1)$$

的前 $m+1$ 行和 $m+1$ 列与 $P_{m+1} \cdots P_1 A_1$ 的前 $m+1$ 行列是不变

化的，因此

$$P_{m+1} \cdots P_1 A_1 P_1$$

的第 1 列就是最后输出的第 1 列，在这之后， P_1 的信息就不再有用而可以退出存区了，因此，计算如果按

$$\cdots (P_{m+1} (P_{m+1} \cdots P_1 A_1) P_1) P_2 \cdots \quad (3.5.34)$$

的次序进行，就可以不必把 P_1, \cdots, P_{n-1} 的全部信息存贮起来，而只需准备一个存贮

$$P_1, P_2, \cdots, P_{m+1}$$

的存区，记为 \overline{W} 存区，其容量为 $(m+1)^2$ 。此外，还应有对

$$P_{m+1} \cdots P_{m+2} (P_{m+1} \cdots P_1 A) P_1 \cdots P_i$$

进行加工的辅助存区。其一，实现由 A_1 计算 R ，称为辅助存区 I；其二，实现由 R 计算 A_2 ，称为辅助存区 II，二者的容量均为 $(m+1) \times (2m+1)$ 。

由 A_1 计算 A_2 的主要步骤如下：

算法 3.5.1

1) 辅区 I 清零，辅区 II 清零；把对称带状阵 A 的上三角部分元素 $(i, i+k)$ 依压缩形式放入二维数组 $B[1:n, 0:m]$ 的 $B[i, k]$ 中， $(i=1, 2, \cdots, n, k=0, 1, \cdots, \min(n-i, m))$ ，以 $n=6, m=2$ 为例，则 A, B 的初态各为

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & & & \\ a_{12} & a_{22} & a_{23} & a_{24} & & \\ a_{13} & a_{23} & a_{33} & a_{34} & a_{35} & \\ & a_{24} & a_{34} & a_{44} & a_{45} & a_{46} \\ & & a_{35} & a_{45} & a_{55} & a_{56} \\ & & & a_{46} & a_{56} & a_{66} \end{pmatrix},$$

$$B = \begin{bmatrix} b_{10} & b_{11} & b_{12} \\ b_{20} & b_{21} & b_{22} \\ b_{30} & b_{31} & b_{32} \\ b_{40} & b_{41} & b_{42} \\ b_{50} & b_{51} & 0 \\ b_{60} & 0 & 0 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{22} & a_{23} & a_{24} \\ a_{33} & a_{34} & a_{35} \\ a_{44} & a_{45} & a_{46} \\ a_{55} & a_{56} & 0 \\ a_{66} & 0 & 0 \end{bmatrix},$$

2) 确定位移量 k_i 。欲求 A 的依模最小的特征值时, 可取 $k_i = 0$; 欲求 A 的在某 t 值附近的特征值时, 可取 $k_i = t$ 。并置 $b_{i0} \leftarrow (b_{i0} - k_i)$ ($i = 1, \dots, n$);

3) 为辅区 I 准备初态。即把 B 的前 $m+1$ 行抄入辅区 I 的对角线右侧的各行, 并使下三角部分与上三角部分对称。以 $m=2$ 为例, 此时辅区 I 的初态为

$$\begin{array}{c} \text{辅 助 存 区 I (初态)} \\ \left(\begin{array}{ccccc} b_{10} - k_i & b_{11} & b_{12} & 0 & 0 \\ & b_{11} & b_{20} - k_i & b_{21} & b_{22} & 0 \\ & & b_{12} & b_{21} & b_{30} - k_i & b_{31} & b_{32} \end{array} \right) \end{array}$$

以下按 (3.5.34) 的计算次序实现算法。

4) 对 $i = 1, 2, \dots, m+n$ 做下列循环:

当 $i \leq n$ 时做 4.1), 否则转 4.4);

4.1) 用辅区 I 的第 1 列做单位向量 w_i 及矩阵 $P = I - 2w_i w_i^T$, 用 P 左乘辅区 I;

4.2) 当 $i \leq m+1$ 时做 4.3), 否则转 4.4);

4.3) 送 w_i 到 W 存区的第 i 行; 辅区 I 的第 1 行送入辅区 II 的第 i 行对角线右侧; 转 4);

4.4) 用 W 的第 1 行 w 做 $P = I - 2w w^T$, 用 P 右乘辅区 II, 送辅区 II 第 1 列到 B 的 $i-m-1$ 行, 辅区 II 上移 1 行, 左移 1 列。末行清零, 辅区 I 的第 1 行前 $m+1$ 个元素送入辅区 II 末

行对角线右侧，辅区 I 上移 1 行，左移 1 列，末行清零；W 区上移一行；

4.5) 当 $i < n$ 时，做 4.6)，否则转 4)

4.6) w_i 送入 W 区末行；B 的第 i 行送辅区 I 末行对角线右侧，B 的第 i 条对角线（走向向上），即 $b_{i,1}, b_{i-1,2}, \dots, b_{i-m+1,m}$ 分别送入辅区 I 中 $(m+1, m), (m+1, m-1), \dots, (m+1, 1)$ 的位置；转 4)。

以上步骤完成时，在数组 B 中得到的将是 A 的带位移的 QR 分解的压缩存贮形式。

第三章 习 题

3.1 将矩阵

$$A = \begin{pmatrix} 0 & 1 & & \\ 1 & 0 & 0.01 & \\ & 0.01 & 0 & 1 \\ & & 1 & 0 \end{pmatrix}, \quad \begin{pmatrix} 1 & 0.01 & & \\ 0.01 & 1 & 0.01 & \\ & 0.01 & 1 & 0.01 \\ & & 0.01 & 1 \end{pmatrix}$$

各做两次带位移的 QR 分解，位移量取为矩阵右下角的元素。

3.2 试证，矩阵

$$P = \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}$$

的 QR 分解序列是不变的。

3.3 求矩阵

$$\Lambda = \begin{pmatrix} & & \lambda_1 \\ & \lambda_2 & \\ \lambda_3 & & \end{pmatrix}$$

的 QR 分解序列，它收敛否？

3.4 不可约的 Hessenberg 阵是非减次的，证明之。

3.5 构造矩阵 A 使其 LR 分解是发散的.

3.6 证明三对角对称阵

$$S = \begin{pmatrix} d_1 & e_2 & & \\ e_2 & d_2 & e_3 & \\ & e_3 & \ddots & \ddots \\ & & \ddots & \ddots & e_n \\ & & & e_n & d_n \end{pmatrix}$$

的所有元素满足, $|d_i|, |e_i| \leq \max |\lambda_i|$, 这里 λ_i 是矩阵 S 的第 i 个特征值. 又 S 的 QR 分解中的上三角阵 R 的元素满足 $|r_{ij}| \leq 2 \max |\lambda_i|$.

3.7 设 B 是实对称矩阵, $B = QR$ 是其 QR 分解.

试证

$$|r_{ii}| \geq \min |\lambda_i(B)|.$$

3.8 设 A 是不可约的奇异的 Hessenberg 阵, A_2 是 A 的 QR 分解的第 2 项.

试证 A_2 的最后一行是零行. 设 A 是不可约的且以 $\lambda = 0$ 为 m 重特征值的 Hessenberg 阵, 试推广所证的结果.

3.9 设 $\det(A) \neq 0$, 则存在排列阵 P 使 PA 有 LR 分解, 且其 LR 分解中的单位下三角阵的元素依模不超过 1.

3.10 设 $A_i \rightarrow I$, $A_i = Q_i R_i$ 是 A_i 的 QR 分解. 证明, $Q_i \rightarrow I$, $R_i \rightarrow I$.

3.11 依公式 (3.4.3) ~ (3.4.9) 编制一个求上 Hessenberg 阵的带两步位移的 QR 分解程序. 并用它计算

$$A_1 = \begin{pmatrix} 5 & -2 & -5 & -1 \\ 1 & 0 & -3 & 2 \\ 0 & 2 & 2 & -3 \\ 0 & 0 & 1 & -2 \end{pmatrix}$$

的 A_3 .

3.12 编制一个求带行交换的计算上 Hessenberg 阵的行列式的程序, 并用它计算

$$\begin{pmatrix} 1 & 0.25 & & \\ 2 & 1 & 0.25 & \\ & 2 & 1 & 0.25 \\ & & 2 & 1 \end{pmatrix}$$

的行列式。

3.13 依本章第5节所述的方法编制一个求对称带状阵的全部特征值的程序，并用它计算

$$\begin{pmatrix} 4 & 2 & 1 & & \\ & 2 & 4 & 2 & 1 \\ & 1 & 2 & 4 & 2 & 1 \\ & & 1 & 2 & 4 & 2 \\ & & & 1 & 2 & 4 \end{pmatrix}$$

的全部特征值。

3.14 设 A 是正定三对角对称阵， A_s 是 A 的 QR 分解序列的第 s 项，设记 A_s 的对角元素为：

$$\alpha_{s,1}, \alpha_{s,2}, \dots, \alpha_{s,n}$$

而次对角元记为：

$$\beta_{s,1}, \beta_{s,2}, \dots, \beta_{s,n-1}.$$

试证， $\alpha_{s+1,1}, \dots, \alpha_{s+n,n}; \beta_{s+1,1}, \dots, \beta_{s+1,n-1}$ 可由 $\alpha_{s,1}, \dots, \alpha_{s,n}, \beta_{s,1}, \dots, \beta_{s,n}$ 经有限次有理运算得到。

第四章 广义特征值问题

§1 引言

设 A 和 B 均为 n 阶矩阵，广义特征值问题就是求满足矩阵方程

$$Ax = \lambda Bx \quad (4.1.1)$$

的数 λ 及非零向量 x ，它们分别称为广义特征值问题 $A - \lambda B$ 的特征值和特征向量。当然应该把 $N(A) \cap N(B)$ 中的非零向量（如果存在的话）排除在外。因为对于这些向量，任意数 λ 都是特征值，显然这是毫无意义的。排除了这种情形，特征值 λ 为特征多项式

$$p_n(\lambda) = \det(A - \lambda B) \quad (4.1.2)$$

的零点。

现在，我们在前章的基础上，对广义特征值问题的计算方法作初步分析：考虑一下解决这一问题有些什么困难与矛盾；设想一下解决这一问题有哪些途径。

最为自然的想法是：如果 B 为非异，那末方程 (4.1.1) 就等价于方程

$$B^{-1}Ax = \lambda x. \quad (4.1.3)$$

这样问题 (4.1.1) 就化为 $B^{-1}A$ 的标准特征值问题。只要 B 关于求逆运算并非病态，那末就可利用任何一种标准特征值问题的计算方法来求解问题 (4.1.3)。

在工程计算中，常常会碰到 A 和 B 均为对称正定的情形（或仅其中之一为对称正定）。譬如，在结构系统的振动问题中， A

为刚度矩阵, B 为质量矩阵, $\lambda = \omega^2$ (ω 为自振频率) 就属于这种情形。这时问题 (4.1.1) 就称为对称广义特征值问题。由于 B 是对称正定阵, 所以对它可以作 Cholesky 分解:

$$B = LL^T$$

这里 L 是对角元皆为正的下三角阵。众所周知, 这一分解是十分稳定的, 于是 (4.1.1) 就等价于

$$L^{-1}AL^{-T}(L^Tx) = \lambda(L^Tx).$$

若令

$$y = L^Tx, \quad (4.1.4)$$

$$C = L^{-1}AL^{-T}, \quad (4.1.5)$$

则 (4.1.1) 问题即等价于标准对称特征值问题:

$$Cy = \lambda y. \quad (4.1.6)$$

这样前面讲的关于对称特征值问题的计算方法都可用来求解 (4.1.6)。例如, 可用 **TRED2** 将 C 化为对称三对角形, 然后利用 **TQL2** 求出 C 的全部特征值和特征向量, 最后利用 **REBAK** 求出原问题的特征向量。

然而, 在以上讨论的各种计算方案中, 蕴含着许多矛盾。首先, 如果 B 对于求逆运算为病态, 那末此时就不能较为精确地计算 $B^{-1}A$ 。从而必然影响特征值、特征向量的精确计算。事实上, 当 B 关于求逆运算为病态时, 一个本质的困难在于: 矩阵 B 的最小奇异值 σ 相对而言较小, 因而问题 (4.1.3) 必有一些相对而言较大的特征值, 而它们对于 B 的扰动是十分敏感的。例如, 设

$$A = \begin{pmatrix} 1 & 2 \\ 2 & 0 \end{pmatrix}, \quad B = \begin{pmatrix} 1+\varepsilon & 1 \\ 1 & 1 \end{pmatrix},$$

其中 $\varepsilon > 0$, 则

$$\det(A - \lambda B) = \varepsilon\lambda^2 + 3\lambda - 4,$$

故 $A - \lambda B$ 的特征值是

$$\lambda_{1,2} = \frac{-3 \pm \sqrt{9 + 16\varepsilon}}{2\varepsilon}.$$

当 $\varepsilon \rightarrow 0$ 时, $A - \lambda B$ 的特征值趋于 $-\frac{4}{3}$ 和 ∞ . 这说明, ε 很小的变化, 即对 B 很小的扰动, 就能引起第二个特征值很大的改变.

为了解决这一矛盾, 显然应避免求逆运算. QZ 算法就是基于这一想法产生的一种算法.

前面, 已介绍了求解标准特征值问题的 QR 算法. 其基本思想是: 采用一系列 (无穷个) 酉相似变换将原始问题化为上三角形或拟上三角形. 由此我们得到启发, 是否能将这种想法移植过来, 以求解广义特征值问题.

因为对 A 和 B 进行相同的酉等价变换并不改变原始问题的特征值. 所以人们想到利用酉等价变换同时将 A, B 化为上三角形 \tilde{A} 和 \tilde{B} . 这样, 问题 (4.1.1) 的特征值即可由 \tilde{A}, \tilde{B} 对角元之比给出, 而整个计算过程并不涉及求逆运算. 但这一算法也有缺点, 它在计算过程中要破坏原始矩阵的性质, 如对称性、稀疏性等等.

对于对称的广义特征值问题, 固然可以化为对称的标准特征值问题. 可是, 实际工程问题中的矩阵 A 和 B 往往是对称带形矩阵, 带宽常常很小. 而将 (4.1.1) 化为 (4.1.6) 意味着带形结构的破坏. 因为一般而言矩阵 C 将为一满阵. 不充分利用 A, B 稀疏带状的特点, 而间接求解 (4.1.6), 这样的方案显然是低效的. 针对 A, B 为稀疏带状的性质, Peters 和 Wilkinson 提出了直接求解对称广义特征值问题的方法. 其基本思想与标准特征值问题的 Givens-Householder 方法是类同的. 简言之, 即利用 Sturm 序列:

$$\begin{aligned} P_0(\lambda) &\equiv 1, \\ P_r(\lambda) &\equiv \det(A^{[r]} - \lambda B^{[r]}) \end{aligned} \quad (4.1.7)$$

的同符号数，结合二分法求得满足给定转度的特征值。而确定 Sturm 序列 (4.1.7) 的同符号数时，可以充分利用 A 、 B 的稀疏带状性质，(4.1.7) 式中 $A^{[r]}$ 表示矩阵 A 的 r 阶首主子矩阵。

在工程计算中，对于振动问题，往往只需求 $A - \lambda B$ 最小的几个特征值及其相应的特征向量（相当于求系统的几个最低的固有频率及其振型）。这时对 (4.1.6) 采用子空间迭代法也是有利的。采用子空间迭代法求解 (4.1.6) 时， $C = L^{-1}AL^{-T}$ 不必真正形成，而只要存贮其因子 L 和矩阵 A 。因为在子空间迭代法中，只用到矩阵 A 与向量的乘积，以及求解系数矩阵为 L 和 L^T 的方程组。其次，它又可以利用矩阵 A 的稀疏性，以提高计算效率。

§2 QZ 算法

按前节的设想，QZ 算法可按以下三个步骤进行：

(1) 首先利用正交三角化算法（例如可采用 Householder 方法或修正的 Gram-Schmidt 算法）将 B 化为上三角形 \tilde{B} ，同时将这些变换施行于 A 化为 \tilde{A} 。如 $n=6$ 时有

$$\begin{array}{cccccc} \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \end{array} \quad \begin{array}{cccccc} \times & \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & 0 & \times; \end{array}$$

\tilde{A}
 \tilde{B}

(2) 利用平面旋转阵逐步把矩阵 \tilde{A} 化为上 Hessenberg 阵，但要求这些平面旋转阵同时作用 \tilde{B} 后仍保持为上三角阵。消去 \tilde{A}

次对角线下的元素，其顺序为：\$(n, 1), (n-1, 1), \dots, (3, 1); (n, 2), (n-1, 2), \dots, (4, 2); \dots, (n, n-2)\$ 一元，消去 \$\bar{A}\$ 的 \$(i, i)\$ 一元，可通过左乘平面旋转阵

$$Q_{i-1,i} = \begin{pmatrix} 1 & & & & 0 \\ & \ddots & & & \\ & & 1 & & \\ & & & \cos\varphi_i & \sin\varphi_i \\ & & & -\sin\varphi_i & \cos\varphi_i \\ & & 0 & & 1 & \ddots & \\ & & & & & & 1 \end{pmatrix} \quad \begin{matrix} i-1 & (4.2.1) \\ i \end{matrix}$$

来实现。而 \$Q_{i-1,i}\$ 作用于 \$\tilde{B}\$ 时，在 \$\tilde{B}\$ 的 \$(i, i-1)\$ 位置上（即次对角线上）引入了一个非零元。但这个非零元可以通过右乘另一平面旋转阵

$$Z_{i-1,i} = \begin{pmatrix} 1 & & & & 0 \\ & \ddots & & & \\ & & 1 & & \\ & & & \cos\theta_i & \sin\theta_i \\ & & & -\sin\theta_i & \cos\theta_i \\ & & 0 & & 1 & \ddots & \\ & & & & & & 1 \end{pmatrix} \quad \begin{matrix} i-1 & (4.2.2) \\ i \end{matrix}$$

消去。然而这个矩阵右乘 \$Q_{i-1,i}\bar{A}\$ 并不破坏已经得到的零元素。以上例来说明，第一步变换过程如下：

$$\begin{array}{cccccccccccc} \times & \times & \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & 0 & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & 0 & 0 & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & 0 & 0 & 0 & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & \times & \times \\ 0 & \times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & \times & \times \end{array} \quad \begin{matrix} Q_{56}\bar{A}Z_{56} & Q_{56}\tilde{B} \end{matrix}$$

$$\begin{array}{cccccc}
\times & \times & \times & \times & \times & \times \\
0 & \times & \times & \times & \times & \times \\
0 & 0 & \times & \times & \times & \times \\
0 & 0 & 0 & \times & \times & \times \\
0 & 0 & 0 & 0 & \times & \times \\
0 & 0 & 0 & 0 & 0 & \times
\end{array}$$

$$Q_{56}\tilde{B}Z_{56}$$

第二步变换过程如下:

$$\begin{array}{cccccc}
\times & \times & \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\
\times & \times & \times & \times & \times & \times & 0 & \times & \times & \times & \times & \times \\
\times & \times & \times & \times & \times & \times & 0 & 0 & \times & \times & \times & \times \\
\times & \times & \times & \times & \times & \times & 0 & 0 & 0 & \times & \times & \times \\
0 & \times & \times & \times & \times & \times & 0 & 0 & 0 & \times & \times & \times \\
0 & \times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & 0 & \times
\end{array}$$

$$Q_{45}Q_{56}\tilde{A}Z_{56}$$

$$Q_{45}Q_{56}\tilde{B}Z_{56}$$

$$\begin{array}{cccccc}
\times & \times & \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\
\times & \times & \times & \times & \times & \times & 0 & \times & \times & \times & \times & \times \\
\times & \times & \times & \times & \times & \times & 0 & 0 & \times & \times & \times & \times \\
\times & \times & \times & \times & \times & \times & 0 & 0 & 0 & \times & \times & \times \\
0 & \times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & \times & \times \\
0 & \times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & 0 & \times
\end{array}$$

$$Q_{45}Q_{56}\tilde{A}Z_{56}Z_{45}$$

$$Q_{45}Q_{56}\tilde{B}Z_{56}Z_{45}$$

如此继续执行以上步骤,最后 \tilde{A} 化为上 Hessenberg 阵而 \tilde{B} 仍保持上三角形。我们把这两个矩阵仍分别记作 \tilde{A} , \tilde{B} 。

(3) 因为

$$\tilde{A} = UAV,$$

$$\tilde{B} = UB\tilde{V},$$

这里 $U, V \in \mathbb{R}^{n \times n}$ 为第一类正交阵. 因为

$$\begin{aligned}\det(\tilde{A} - \lambda \tilde{B}) &= \det(U) \cdot \det(V) \cdot \det(A - \lambda B) \\ &= \det(A - \lambda B).\end{aligned}$$

由此可见, 问题 (4.1.1) 应与问题

$$\tilde{A}x = \lambda \tilde{B}x \quad (4.2.3)$$

有相同的特征值. 另外, 若 x 为问题 (4.2.3) 的特征向量, 则 Vx 为问题 (4.1.1) 对应的特征向量. 因此, 现在只要讨论广义特征值问题 (4.2.3) 的计算方法就行了.

现设 \tilde{B} 为非异, 而 $\tilde{C} \equiv \tilde{A}\tilde{B}^{-1}$. 由于 \tilde{A} 为上 Hessenberg 阵而 \tilde{B} 为上三角形, 故 \tilde{C} 也为上 Hessenberg 阵. 根据 §1 的分析, 现在的目标是要把 \tilde{A} 与 \tilde{B} 同时化约为上三角阵, 也就等价于将 \tilde{C} 化约为上三角阵.

鉴于以上分析, 我们想到利用带原点位移的 QR 算法. 若将它施于 \tilde{C} , 得到矩阵序列 $\{C_k\}$, 则在一定条件下, 该矩阵序列应趋于一上三角阵. 具体地说, 即

$$\begin{cases} \tilde{C} \equiv C_1, \\ Q_k(C_k - \kappa_k I) \equiv R_k, \quad (k = 1, 2, \dots) \\ R_k Q_k^H + \kappa_k I \equiv C_{k+1}, \end{cases} \quad (4.2.4)$$

且有

$$C_{k+1} = Q_k C_k Q_k^H, \quad (4.2.5)$$

然而我们的目的并不是要从 C_k 直接产生 Q_k , 再进而求出 C_{k+1} .

假定 A_k 与 B_k 满足 $A_k B_k^{-1} \equiv C_k$, 且 A_k 与 B_k 仍分别为上 Hessenberg 阵与上三角阵, 那末可以设想这样的计算方案: 避开直接计算 B_k^{-1} , 间接地求出 Q_k , 继而构造 A_{k+1} 和 B_{k+1} :

$$A_{k+1} = Q_k A_k V_k, \quad (4.2.6)$$

$$B_{k+1} = Q_k B_k V_k.$$

这里 V_k 为酉阵，它要保证 A_{k+1} 为上 Hessenberg 阵而 B_{k+1} 为上三角阵。于是

$$\begin{aligned} C_{k+1} &\equiv A_{k+1} B_{k+1}^{-1} \\ &= (Q_k A_k V_k) (Q_k B_k V_k)^{-1} \\ &= Q_k C_k Q_k^H. \end{aligned} \quad (4.2.7)$$

如何避开直接计算 B_k^{-1} 而间接地求出 Q_k 呢？由定理 3.3.1 知：对 (4.2.7) 而言，如果 C_{k+1} 为不可约的上 Hessenberg 阵，则 C_{k+1} 与 Q_k 唯一地（除模为 1 的常数因子外）由 Q 的第一行决定。这样我们就很自然地想到隐位移 QR 算法，为避免复数运算又可采用二步 QR 算法。

据前章分析知，为实现二步 QR 算法，首先应找到一酉阵 P_0 ，它的第一行应与矩阵

$$\varphi_2(C_k) \equiv (C_k - \lambda_k I)(C_k - \lambda_{k+1} I) \quad (4.2.8)$$

的第一列成比例，其中 $\lambda_{k+1} = \bar{\lambda}_k$ 。而

$$\begin{aligned} \varphi_2(C_k) e_1 &= \begin{pmatrix} (\gamma_{11}^{(k)})^2 + \gamma_{21}^{(k)} \gamma_{12}^{(k)} \\ \gamma_{21}^{(k)} \gamma_{11}^{(k)} + \gamma_{22}^{(k)} \gamma_{21}^{(k)} \\ \gamma_{31}^{(k)} & \gamma_{21}^{(k)} \\ 0 \\ \vdots \\ 0 \end{pmatrix} - (\lambda_k + \lambda_{k+1}) \begin{pmatrix} \gamma_{11}^{(k)} \\ \gamma_{21}^{(k)} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \\ &\quad + \lambda_k \lambda_{k+1} \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}. \end{aligned} \quad (4.2.9)$$

由此可见， $\varphi_2(C_k) e_1$ 仅用到 C_k 的前二列。但因 B_k 是上三角形，故 B_k^{-1} 亦为上三角形。于是 C_k 的前二列为

$$(a_1^{(k)}, a_2^{(k)}) \begin{pmatrix} \beta_{11}^{(k)} & \beta_{12}^{(k)} \\ 0 & \beta_{22}^{(k)} \end{pmatrix}^{-1}. \quad (4.2.10)$$

这样，为了计算 P_0 ，只需计算 B_k 的二阶主子矩阵的逆即可。显然，这要比计算 B_k^{-1} 的效果好得多。又由前章知，因为 $\varphi_2(C_k)$ 的第一列至多前三个元素不为零，故 P_0 不妨取作如下形状：

$$P_0 = \begin{pmatrix} R_0 & 0 \\ 0 & I_{n-3} \end{pmatrix}, \quad (4.2.11)$$

其中 R_0 是一个三阶的初等镜像变换阵。设 $\varphi_2(C_k)e_1 \equiv \tilde{q}_1$ ，则 P_0^T 的第一列即为 $\tilde{q}_1 = \pm \frac{\tilde{q}_1}{\|\tilde{q}_1\|}$ ，即 P_0 应具有这样的性质： $P_0^T e_1 = \tilde{q}_1$ 或 $\tilde{q}_0 \tilde{q}_1 = \sigma e_1$ ($\sigma = \pm \|\tilde{q}_1\|$)。当然这样的 P_0 是能构造出来的。 $P_0 A_k$ 和 $P_0 B_k$ 应具有如下形状（仍以 $n=6$ 为例）

$$\begin{array}{cccccc} \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times \end{array} \quad \begin{array}{cccccc} \times & \times & \times & \times & \times & \times \\ \otimes & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & 0 & \times \end{array}$$

$P_0 A_k \qquad P_0 B_k$

根据以上设想方案，为了求出 A_{k+2} 和 B_{k+2} ，需找到这样的正交阵 U_k 和 V_k ，使 $U_k P_0 A_k V_k$ 为上Hessenberg阵， $U_k P_0 B_k V_k$ 为上三角阵。同时，还要保证 $U_k P_0 = Q_k'$ 与 P_0 有相同的第一行，这样 Q_k' 与 Q_k 即有相同的第一行。据定理3.3.1知， $Q_k' \equiv Q_k \cdot U_k$ 和 V_k 的形成步骤如下：

- (1) 首先将 $P_0 B_k$ 按下面方式化为上三角阵（以 $n=6$ 为例）
- 选取初等镜像变换阵 Z_0 ，使 $P_0 B_k Z_0$ 中 $(3,1)$ 元和 $(3,2)$

元变为零。然后再用平面旋转阵 $R_{12} = Z'_0$ 使 $P_0 B_k Z_0 Z'_0$ 中的 (2,1) 元变为零。于是 $P_0 A_k Z_0 Z'_0$ 和 $P_0 B_k Z_0 Z'_0$ 就有如下的形状:

$ \begin{array}{cccccc} \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \otimes & \times & \times & \times & \times & \times \\ \otimes & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times \end{array} $	$ \begin{array}{cccccc} \times & \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & 0 & \times \end{array} $
$P_0 A_k Z_0 Z'_0$	$P_0 B_k Z_0 Z'_0$

请注意这里 Z_0 和 Z'_0 的形状为:

$ \begin{array}{cccccc} \times & \times & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{array} $	$ \begin{array}{cccccc} \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{array} $
Z_0	Z'_0

(2) 选取一初等镜象变换阵 P_1 , 使 $P_0 A_k Z_0 Z'_0$ 中打圈位置上的元素变为零, P_1 的形状如下:

$$\begin{array}{cccccc}
 1 & 0 & 0 & 0 & 0 & 0 \\
 0 & \times & \times & \times & 0 & 0 \\
 0 & \times & \times & \times & 0 & 0 \\
 0 & \times & \times & \times & 0 & 0 \\
 0 & 0 & 0 & 0 & 1 & 0 \\
 0 & 0 & 0 & 0 & 0 & 1
 \end{array}$$

P_1

于是 $P_1 P_0 A_k Z_0 Z'_0$ 和 $P_1 P_0 B_k Z_0 Z'_0$ 有如下形状:

$$\begin{array}{c|cccccc}
 & \times & \times & \times & \times & \times & \times \\
 \hline
 \times & \times & \times & \times & \times & \times & \times \\
 0 & \times & \times & \times & \times & \times & \times \\
 0 & \times & \times & \times & \times & \times & \times \\
 0 & 0 & 0 & \times & \times & \times & \times \\
 0 & 0 & 0 & 0 & \times & \times & \times
 \end{array}
 \quad
 \begin{array}{c|cccccc}
 & \times & \times & \times & \times & \times & \times \\
 \hline
 0 & \times & \times & \times & \times & \times & \times \\
 0 & \otimes & \times & \times & \times & \times & \times \\
 0 & \otimes & \otimes & \times & \times & \times & \times \\
 0 & 0 & 0 & 0 & \times & \times & \times \\
 0 & 0 & 0 & 0 & 0 & \times & \times
 \end{array}$$

(3) 对划去第一行、第一列后的子矩阵重复 (1)(2) 步骤。

(4) 当算法进行到第 $n-2(=4)$ 次循环第(1)步时, 二个矩阵呈如下形状:

$$\begin{array}{c|cccc}
 \times & \times & \times & \times & \times \\
 \times & \times & \times & \times & \times \\
 0 & \times & \times & \times & \times \\
 \hline
 0 & 0 & \times & \times & \times \\
 0 & 0 & 0 & \times & \times \\
 0 & 0 & 0 & \otimes & \times
 \end{array}
 \quad
 \begin{array}{c|cccc}
 \times & \times & \times & \times & \times \\
 0 & \times & \times & \times & \times \\
 0 & 0 & \times & \times & \times \\
 \hline
 0 & 0 & 0 & \times & \times \\
 0 & 0 & 0 & 0 & \times \\
 0 & 0 & 0 & 0 & 0
 \end{array},$$

选取一平面旋转阵 P_4 , 使左边矩阵中打圈位置上的元素变为零。 P_4 的形状如下:

$$\begin{array}{cccccc}
 1 & 0 & 0 & 0 & 0 & 0 \\
 0 & 1 & 0 & 0 & 0 & 0 \\
 0 & 0 & 1 & 0 & 0 & 0 \\
 0 & 0 & 0 & 1 & 0 & 0 \\
 0 & 0 & 0 & 0 & \times & \times \\
 0 & 0 & 0 & 0 & \times & \times
 \end{array}$$

P_k 左乘这二个矩阵后有如下形状:

$$\begin{array}{ccc|ccc}
 \times & \times & \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\
 \times & \times & \times & \times & \times & \times & 0 & \times & \times & \times & \times & \times \\
 0 & \times & \times & \times & \times & \times & 0 & 0 & \times & \times & \times & \times \\
 \hline
 0 & 0 & \times & \times & \times & \times & 0 & 0 & 0 & \times & \times & \times \\
 0 & 0 & 0 & \times & \times & \times & 0 & 0 & 0 & 0 & \times & \times \\
 0 & 0 & 0 & 0 & \times & \times & 0 & 0 & 0 & 0 & \otimes & \times
 \end{array}$$

(5) 选取一平面旋转阵 Z_k , 右乘右边那个矩阵, 使其中打圈位置上的元素变为零. 同时 Z_k 右乘左边的矩阵, 显然它仍保持上 **Hessenberg** 阵的形式. 至此化约完成, 即 A_{k+2} 和 B_{k+2} 均已求得.

以上全部过程构成了所谓二步 QZ 算法. 实际上, 它是二步 QR 算法的一种推广.

现在进一步讨论计算细节. 根据以上分析. 我们知道, 对 (4.2.3) 施行二步 QZ 算法, 首先应找到正交阵 P_0 , 它的第一行应与矩阵 $\Phi_2(C_k)$ 的第一列成比例. 为此应找出 C_k 的前二列. 据 (4.2.10) 知, C_k 的前二列应为:

$$\begin{pmatrix}
 a_{11}^{(k)}/\beta_{11}^{(k)} & a_{12}^{(k)}/\beta_{22}^{(k)} - a_{11}^{(k)}\beta_{11}^{(k)}/\beta_{11}^{(k)}\beta_{22}^{(k)} \\
 a_{21}^{(k)}/\beta_{11}^{(k)} & a_{22}^{(k)}/\beta_{22}^{(k)} - a_{21}^{(k)}\beta_{12}^{(k)}/\beta_{11}^{(k)}\beta_{22}^{(k)} \\
 \vdots & \vdots \\
 0 & 0
 \end{pmatrix} \quad (4.2.12)$$

为了确定位移量 κ_k 和 κ_{k+1} , 则应确定矩阵 C_k 的尾二阶主子矩阵, 我们把它近似地取作:

$$\begin{pmatrix}
 a_{n-1,n-1}^{(k)} & a_{n-1,n}^{(k)} \\
 a_{n,n-1}^{(k)} & a_{n,n}^{(k)}
 \end{pmatrix}
 \begin{pmatrix}
 \beta_{n-1,n-1}^{(k)} & \beta_{n-1,n}^{(k)} \\
 0 & \beta_{n,n}^{(k)}
 \end{pmatrix}^{-1}. \quad (4.2.13)$$

位移量 κ_k 和 $\kappa_{k+1} \cdots$ 为矩阵 (4.2.13) 的二个特征值。据二次方程根与系数的关系即得:

$$\begin{cases} \kappa_k + \kappa_{k+1} = \alpha_{nn}^{(k)} / \beta_{nn}^{(k)} + \alpha_{n-1,n-1}^{(k)} / \beta_{n-1,n-1}^{(k)} \\ \quad - \alpha_{n,n-1}^{(k)} \beta_{n-1,n}^{(k)} / (\beta_{n-1,n-1}^{(k)} \beta_{nn}^{(k)}), \\ \kappa_k \cdot \kappa_{k+1} = (\alpha_{n-1,n-1}^{(k)} \alpha_{nn}^{(k)} - \alpha_{n,n-1}^{(k)} \alpha_{n-1,n}^{(k)}) / (\beta_{n-1,n-1}^{(k)} \beta_{nn}^{(k)}). \end{cases} \quad (4.2.14)$$

将 (4.2.12) 和 (4.2.14) 代入 (4.2.9) 即得 \tilde{q}_1 前三个非零分量:

$$\begin{cases} \alpha_{10} = \left[\left(\frac{\alpha_{n-1,n-1}^{(k)}}{\beta_{n-1,n-1}^{(k)}} - \frac{\alpha_{11}^{(k)}}{\beta_{11}^{(k)}} \right) \left(\frac{\alpha_{nn}^{(k)}}{\beta_{nn}^{(k)}} - \frac{\alpha_{11}^{(k)}}{\beta_{11}^{(k)}} \right) - \frac{\alpha_{n,n-1}^{(k)} \alpha_{n-1,n}^{(k)}}{\beta_{nn}^{(k)} \beta_{n-1,n-1}^{(k)}} \right. \\ \quad \left. + \frac{\alpha_{n,n-1}^{(k)} \beta_{n-1,n}^{(k)}}{\beta_{nn}^{(k)} \beta_{n-1,n-1}^{(k)}} \cdot \frac{\alpha_{11}^{(k)}}{\beta_{11}^{(k)}} \right] \frac{\beta_{11}^{(k)}}{\alpha_{21}^{(k)}} + \frac{\alpha_{12}^{(k)}}{\beta_{22}^{(k)}} - \frac{\alpha_{11}^{(k)} \beta_{22}^{(k)}}{\beta_{11}^{(k)} \beta_{22}^{(k)}}, \\ \alpha_{20} = \frac{\alpha_{22}^{(k)}}{\beta_{22}^{(k)}} + \frac{\alpha_{11}^{(k)}}{\beta_{11}^{(k)}} - \frac{\alpha_{21}^{(k)} \beta_{12}^{(k)}}{\beta_{11}^{(k)} \beta_{22}^{(k)}} - \frac{\alpha_{nn}^{(k)}}{\beta_{nn}^{(k)}} - \frac{\alpha_{n-1,n-1}^{(k)}}{\beta_{n-1,n-1}^{(k)}} \\ \quad + \frac{\alpha_{n,n-1}^{(k)} \beta_{n-1,n}^{(k)}}{\beta_{nn}^{(k)} \beta_{n-1,n-1}^{(k)}}, \\ \alpha_{30} = \frac{\alpha_{32}^{(k)}}{\beta_{22}^{(k)}}. \end{cases} \quad (4.2.15)$$

注意, 首先, \tilde{q}_1 的每个分量我们都乘以 $\beta_{11}^{(k)} / \alpha_{21}^{(k)}$, 这对确定 P_0 并无影响. 确定 P_0 的具体过程, 见前章 §3; 其次, 为了程序设计上的原因, 把 \tilde{q}_1 的三个非零分量分别记作 $\alpha_{10}, \alpha_{20}, \alpha_{30}$, 把 $P_0 A_k$ 第一列的元素记作 α_{i1} , 把 $P_1 P_0 A_k Z_0 Z'_0 Z_1 Z'_1$ 的第二列记作 α_{i2} 等等. 至此, 我们即可构成二步 QZ 算法.

算法 4.2.1 设给定矩阵 $A, B \in R^{n \times n}$ 而 B 为非异矩阵, 用 QZ 算法求 (4.1.1) 的全部特征值. 用此算法求解的步骤如下:

1) 利用正交三角化算法将 B 化为上三角阵 \tilde{B} , 同时将 A 化为 \tilde{A} . 置 $A = \tilde{A}, B = \tilde{B}$;

2) 利用平面旋转阵将 A 化为上 Hessenberg 阵, 同时保持 B 为上三角形;

3) 对 $j = n, n-1, \dots, 1$.

3.1) 对 $k = 1, 2, \dots, m$ (允许的最大迭代次数).

3.1.1) 按 (4.2.15) 计算 $\alpha_{10}, \alpha_{20}, \alpha_{30}$.

3.1.2) 对 $i = 1, 2, \dots, n-2$;

3.1.2.1) 确定 P_{i-1} 消去 $\alpha_{i+1,i-1}^{(k)}, \alpha_{i+2,i-1}^{(k)}$, 形成相应的矩阵.

3.1.2.2) 确定 Z_{i-1} 消去 $\beta_{i+2,i+1}^{(k)}, \beta_{i+2,i}^{(k)}$, 形成相应的矩阵.

3.1.2.3) 确定 Z'_{i-1} 消去 $\beta_{i+1,i}^{(k)}$, 形成相应的矩阵.

3.1.3) NEXT i .

3.1.4) 确定 P_{n-2} 消去 $\alpha_{n,n-2}^{(k)}$, 形成相应的矩阵.

3.1.5) 确定 Z_{n-2} 消去 $\beta_{n,n-1}^{(k)}$, 形成相应的矩阵.

3.1.6) 如果满足收敛准则, 那末置 $\lambda_i = \alpha_{ii}^{(k)} / \beta_{ii}^{(k)}$ 并转 4).

3.2) NEXT k

4) NEXT j

5) END.

以上算法仅是一个简单的示意描述. 许多具体细节, 读者可进一步参阅 Moler, C. B. 和 Stewart G. W. 的文章 [9]. 这里还应特别指出的是, 和 QR 算法类似, 在 QZ 迭代过程中, 每迭代一次应由下到上地检查 A_k 的次对角元素. 若某次对角元可忽略, 则可将问题分裂为二个低阶的问题, 从而分别进行迭代.

对于 (4.2.3), 其特征向量容易算出. 事实上, 经过有限次酉等价变换后迭代终止, 则有

$$\tilde{A} = \tilde{U} \tilde{A} \tilde{V}, \quad (4.2.16)$$

$$\tilde{B} = \tilde{U} \tilde{B} \tilde{V}. \quad (4.2.17)$$

这里 \tilde{A} 与 \tilde{B} 可以看作上三角阵 (\tilde{A} 的次对角元相当小), 而 \tilde{U} 与 \tilde{V} 为正交阵. 据 (4.2.16) 和 (4.2.17) 有

$$\tilde{B}^{-1} \tilde{A} \tilde{V} = \tilde{V} \tilde{B}^{-1} \tilde{A}. \quad (4.2.18)$$

若设 $\tilde{B}^{-1} \tilde{A} \tilde{z}_i = \lambda_i \tilde{z}_i$, 这里 λ_i 为问题 (4.1.1) 的广义特征值, 那末据 (4.2.18) 即有

$$\begin{aligned} \tilde{B}^{-1} \tilde{A} (\tilde{V} \tilde{z}_i) &= \tilde{V} \tilde{B}^{-1} \tilde{A} \tilde{z}_i \\ &= \lambda_i (\tilde{V} \tilde{z}_i). \end{aligned} \quad (4.2.19)$$

由此可见, $\tilde{V} \tilde{z}_i$ 即为问题 (4.2.3) 对于 λ_i 的广义特征向量. 而 $\tilde{B}^{-1} \tilde{A}$ 的特征向量可由齐次方程

$$(\tilde{A} - \lambda_i \tilde{B}) \tilde{z}_i = 0 \quad (4.2.20)$$

确定之. 综上所述, 问题 (4.1.1) 对应 λ_i 的广义特征向量为 $\tilde{V} \tilde{z}_i$. 顺便指出, QZ 算法需二个 $n \times n$ 辅助数组, 因此当 n 很大时, 存贮量较大.

§3 Peters-Wilkinson 方法

正如引言中指出那样, 为了不破坏矩阵 A 、 B 的对称带状结构, Peters 和 Wilkinson 将标准对称特征值问题的 Givens-Householder 方法的思想用来解决对称广义特征值问题. 为此, 我们首先回顾一下标准对称特征值问题的 Givens-Householder 方法的基本依据:

(1) 设 A 为对称阵, 它具有相异的特征值, 则其首主子矩阵行列式序列

$$P_0(\lambda), P_1(\lambda), \dots, P_k(\lambda), \dots, P_n(\lambda) \quad (4.3.1)$$

为 Sturm 序列, 其中 $P_0(\lambda) = 1$, $P_k(\lambda) = \det(A^{1:k1} - \lambda I)$;

(2) 在以上的假定下, 记 $a(\sigma)$ 为序列 (4.3.1) 在 σ 处相邻两项同号的数目, 则矩阵 A 恰有 $a(\sigma)$ 个特征值不小于 σ .

现在设法将以上结论过渡到广义特征值问题. 如 §1 中所说, 对称广义特征值问题等价于对称标准特征值问题 (4.1.6). 因此, 如果矩阵 C 也具有相异的特征值, 那末即可利用以上二个结论. 当然, 为了不破坏矩阵 A 、 B 的对称带状结构, 我们并不希望真正去构成矩阵 C , 因此就不能直接利用序列 $\{\det(C^{[k]} - \lambda I)\}$ 的同号数来隔离 $(A - \lambda B)$ 的特征值. 然而, $\det(A^{[k]} - \lambda B^{[k]})$ 与 $\det(C^{[k]} - \lambda I)$ ($k = 1, 2, \dots, n$) 在以上假设条件下, 具有相同的符号, 这样问题也就迎刃而解了. 事实上, 若设 B 的 Cholesky 分解为 $B = LL^T$, 则

$$L^{[k]}L^{[k]T} = B^{[k]},$$

$$(L^{[k]})^{-1}A^{[k]}(L^{[k]})^{-T} = C^{[k]},$$

于是

$$\begin{aligned} \det(A^{[k]} - \lambda B^{[k]}) &= \det[L^{[k]} \{ (L^{[k]})^{-1}A^{[k]}(L^{[k]})^{-T} - \lambda I \} L^{[k]T}] \\ &= \det[L^{[k]}(C^{[k]} - \lambda I)L^{[k]T}] \\ &= [\det(L^{[k]})]^2 \cdot \det(C^{[k]} - \lambda I). \end{aligned}$$

因为 $L^{[k]}$ ($k = 1, \dots, n$) 为非异, 故 $\det(A^{[k]} - \lambda B^{[k]})$ 的符号必与 $\det(C^{[k]} - \lambda I)$ ($k = 1, 2, \dots, n$) 的符号一致. 根据以上事实即可得到如下结论:

在以上假设条件下, 不小于 σ 的 $A - \lambda B$ 特征值的数目恰为序列

$$\det(A^{[k]} - \sigma B^{[k]}) \quad (k = 0, 1, \dots, n) \quad (4.3.2)$$

在 σ 处的同号数 $a(\sigma)$, 这里规定 $\det(A^{[0]} - \lambda B^{[0]}) = 1$.

但在这里应特别指出，序列 (4.3.1) 和序列 (4.3.2) 在计算上是有区别的。前者因为矩阵 A 可用初等镜象变换化约为对称三对角阵，所以序列 (4.3.1) 可以递推计算。可是序列 (4.3.2) 的计算却无此方便。

通常采用列主元素消去法计算 $\det(A^{[k]} - \sigma B^{[k]})$ 的值，以保证计算的稳定性。但应注意，这时列主元素消去法是在 k 阶主子矩阵中进行的，而 k 是从 2 逐次递增至 n 。所以，这里列主元素消去法与求解线性方程组时的情形略有不同。其算法简单描述如下：

算法 4.3.1 设给定对称矩阵 $A, B \in R^{n \times n}$ ， B 为正定，欲求 $\det(A^{[k]} - \sigma B^{[k]})$ 的符号 $s_k (k = 1, 2, \dots, n)$ ，记 $C = A - \sigma B$ ， σ 为给定的数。

1) 置 $s_1 = \text{sign } \gamma_{11}$ 。

2) 对 $k = 1, 2, \dots, n-1$

2.1) 置 $s_{k+1} = s_k$ 。

2.2) 对 $i = 1, 2, \dots, k$

2.2.1) 若 $|\gamma_{k+1,i}| > \gamma_{ii}$ ，则 $\gamma_{k+1,i} \leftrightarrow \gamma_{i,i} (j = i, i+1, \dots, n)$ ，否则转 2.2.3)。

2.2.2) $\gamma_{k+1,i}$ 与 γ_{ii} 同号，则置 $s_{k+1} = -s_{k+1}$ 。

2.2.3) 置 $\gamma_{k+1,i} = \frac{\gamma_{k+1,i}}{\gamma_{ii}}$ ，

2.2.4) 置 $\gamma_{k+1,j} = \gamma_{k+1,j} - \gamma_{k+1,i} \gamma_{i,j} (j = i+1, \dots, n)$ ，

2.3) NEXT i ，

2.4) 如果 $\gamma_{k+1,k+1} < 0$ ，则置 $s_{k+1} = -s_{k+1}$ 。

3) NEXT k 。

4) END。

例 4.3.1 利用算法 4.3.1，确定对称矩阵

$$C = \begin{pmatrix} 1 & 4 & 0 & 0 \\ 4 & 2 & 1 & 0 \\ 0 & 1 & 3 & 4 \\ 0 & 0 & 4 & 3 \end{pmatrix}$$

各阶首主子式的符号.

解 显然 $s_0 = 1$, $s_1 = 1$. 按语句 (2.2.1—2.2.4) 执行如下: $k=1$ 时变换为:

$$C \xrightarrow{\substack{\text{第一行与第} \\ \text{二行交换}}} \begin{pmatrix} 4 & 2 & 1 & 0 \\ 1 & 4 & 0 & 0 \\ 0 & 1 & 3 & 4 \\ 0 & 0 & 4 & 3 \end{pmatrix} \xrightarrow{\substack{\text{第二行减去} \\ \frac{1}{4} \times \text{第一行}}} \begin{pmatrix} 4 & 2 & 1 & 0 \\ 0 & 7/2 & -1/4 & 0 \\ 0 & 1 & 3 & 4 \\ 0 & 0 & 4 & 3 \end{pmatrix}$$

因为进行过一次行交换且 γ_{11} 与 γ_{21} 同号, 同时 $\gamma_{22} > 0$, 故 $s_2 = -1$. $k=2$ 时变换为

$$\xrightarrow{\substack{\text{第三行减去} \\ \frac{2}{7} \times \text{第二行}}} \begin{pmatrix} 4 & 2 & 1 & 0 \\ 0 & 7/2 & -1/4 & 0 \\ 0 & 0 & 43/14 & 4 \\ 0 & 0 & 4 & 3 \end{pmatrix}$$

因为没有进行过行交换, 而 $43/14 > 0$, 故据语句 2.1) 知 $s_3 = -1$. $k=3$ 时变换为

$$\xrightarrow{\substack{\text{第三行与第} \\ \text{四行交换}}} \begin{pmatrix} 4 & 2 & 1 & 0 \\ 0 & 7/2 & -1/4 & 0 \\ 0 & 0 & 4 & 3 \\ 0 & 0 & 43/14 & 4 \end{pmatrix} \xrightarrow{\substack{\text{第四行减} \\ \text{去 } 43/56 \\ \times \text{第三行}}} \begin{pmatrix} 4 & 2 & 1 & 0 \\ 0 & 7/2 & -1/4 & 0 \\ 0 & 0 & 4 & 3 \\ 0 & 0 & 0 & 95/56 \end{pmatrix}$$

因为进行过一次交换, γ_{33} 与 γ_{43} 同号而 $95/56 > 0$, 所以 $s_4 = 1$.

从以上简单的例子可以看出，P-W 算法的主要优点是：能够充分利用原始矩阵 A 和 B 的稀疏带状的特点，从而省去了许多存贮量与计算量。因此，它适用于求解大型对称广义特征值问题。事实上，若记 $A - \sigma B$ 的半带宽为 $m+1$ (包括对角元)，则算法 4.3.1 的第 k 步仅涉及第 $k-m+1$ 到第 $k+1$ 行的带内元素。尽管由于行交换，原来矩阵的半带宽可能要变大，但最大不会超过 $2m+1$ 。

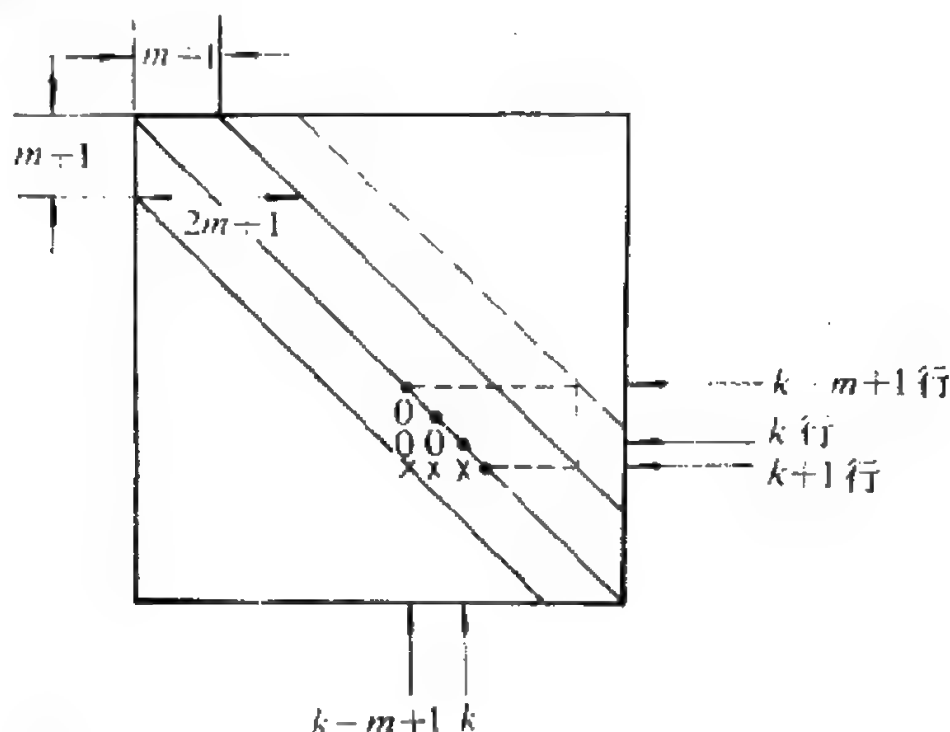


图 4.1 $A - \sigma B$ 在列主元消去过程中带宽变化

P-W 算法的结构基本上和 Givens-Householder 算法是一致的，即利用算法 4.3.1 确定序列 (4.3.2) 的同符号数，然后用二分法来隔离特征值 (部分或全部)。二分法可以与求特征向量的反迭代法以及 Rayleigh 商迭代结合起来，相辅相成。二分法不必进行到隔离出精度范围的特征值为止，而只需进行到将特征值隔离开为止。然后，用反迭代求近似特征向量，再利用 Rayleigh 商使特征值精确化，如此往复循环，直至近似特征值符合精度标准为止。

因为当 B 为非异时, 广义特征值问题(4.1.1)与标准特征值问题(4.1.3)是等价的*. 而后者的反迭代格式是:

$$(B^{-1}A - \tilde{\lambda}_k I) \mathbf{x}_k^{(i+1)} = \mathbf{x}_k^{(i)} \tau_i \quad (i = 0, 1, \dots) \quad (4.3.3)$$

其中 $\sigma_k^{(i)}$ 为位移量, 而 τ_i 为规范化常数. (4.3.3)式两边左乘以 B , 即得

$$(A - \tilde{\lambda}_k B) \mathbf{x}_k^{(i+1)} = B \mathbf{x}_k^{(i)} \tau_i \quad (i = 0, 1, \dots), \quad (4.3.4)$$

这就是广义特征值问题 $A - \lambda B$ 的反迭代格式. 和标准特征值问题一样, 由方程组(4.3.4)求得 $\mathbf{x}_k^{(i+1)}$ 后, 利用Reyleigh商

$$\rho(\mathbf{x}_k^{(i+1)}) = \frac{\mathbf{x}_k^{(i+1)T} A \mathbf{x}_k^{(i+1)}}{\mathbf{x}_k^{(i+1)T} B \mathbf{x}_k^{(i+1)}}, \quad (4.3.5)$$

将特征值精确化最为合适(参阅下节及习题4.6). 根据以上分析即可得到下面的算法.

算法4.3.2 设 $A, B \in \mathbb{R}^{n \times n}$ 且 B 为正定, 又假定广义特征值问题 $A - \lambda B$ 具有相异的特征值, 求 $A - \lambda B$ 第 k 个特征值及其对应的特征向量.

1) 设 σ 为 λ_k 的一个下界, 利用算法4.1.2计算序列(4.3.2)的符号, 以同符号数 $\alpha(\sigma)$ 确定 $[\sigma, \infty)$ 中 $A - \lambda B$ 的特征值的个数(NUMBER子程序);

2) 用二分法将 $A - \lambda B$ 的第 k 个特征值隔离在区间 $[\lambda_k^-, \lambda_k^+]$ 内, 取初始特征值为 $\lambda_k^{(0)} = \frac{1}{2}(\lambda_k^- + \lambda_k^+)$ 或利用割线法确定 $\lambda_k^{(0)}$ (BISECT和INTPOL子程序);

3) 选取初始特征向量 $\mathbf{x}_k^{(0)}$. 通常的选取法是, 用部分主元消去将 $(A - \lambda_k^{(0)} B)$ 化为上三角阵 U , 记下所作的行交换, 然后求解三角形方程组

* 这里当然亦可将对称广义特征值问题化为对称标准特征值问题, 进而考虑它的反迭代格式.

$$U\mathbf{x}_k^{(0)} = \mathbf{e}. \quad (4.3.6)$$

这里 \mathbf{e} 是分量全为 1 的向量;

4) 对 $i = 0, 1, \dots$ (INVERS 子程序),

4.1) 将 $\mathbf{x}_k^{(i)}$ 规范化, 得 $\hat{\mathbf{x}}_k^{(i)}$.

4.2) 利用 $(A - \lambda_k^{(0)}B)$ 的 LU 分解及行交换信息 求解线性方程组

$$(A - \lambda_k^{(0)}B)\mathbf{x}_k^{(i+1)} = B\hat{\mathbf{x}}_k^{(i)},$$

4.3) 利用 Rayleigh 商使 $\lambda_k^{(i)}$ 精确化

$$\lambda_k^{(i+1)} = \frac{\mathbf{x}_k^{(i+1)T} A \mathbf{x}_k^{(i+1)}}{\mathbf{x}_k^{(i+1)T} B \mathbf{x}_k^{(i+1)}}$$

4.4) 收敛判别: 如果

$$\frac{|\lambda_k^{(i+1)} - \lambda_k^{(i)}|}{|\lambda_k^{(i+1)}|} \leq \varepsilon.$$

(这里 ε 是给定的精度) 时, $\lambda_k^{(i+1)}$ 与 $\mathbf{x}_k^{(i+1)}$ 即为所求的结果. 转 6).

5) NEXT i ;

6) 置 $\lambda_k = \lambda_k^{(i+1)}$, $\mathbf{x}_k = \mathbf{x}_k^{(i+1)}$;

7) END.

§4 子空间迭代法

乘幂法和反乘幂法(即反迭代法)是计算标准特征值问题的二个有效算法. 它们特别适应大型稀疏矩阵的特征问题, 但它们只能求得占优特征值和特征向量. 然而实际问题往往需要求前几个占优特征值及其对应的特征向量, 如多元统计分析中的主成分分析等. 过去常采用降阶法, 逐一求出前几个占优特征

值和特征向量。但是，降阶法对各种不同结构的矩阵并无统一的算法；另外，降阶还将破坏原矩阵的形状以及受到舍入误差较大的影响。

子空间迭代法在一定程度上能克服降阶法这些缺点。其基本思想是：对若干个试验向量（可看作 R^n 中某子空间的基）同时进行迭代，从而同时得到前几个占优特征值及其对应的特征向量（或者说 A 的占优子空间）。所以，子空间迭代法可看成幂法、反幂法的一种推广。

这里我们仅讨论实对称阵的情形，至于非对称阵的子空间迭代法，读者可参阅 Stewart, G. W. 以及 Cline, M., Jennings, A. 等的文章[10]、[11]。为说明方便起见，设 $A \in R^{n \times n}$ 为对称正定阵而

$$0 < \lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_n$$

和

$$x_1, x_2, \dots, x_n$$

分别为它的特征值和对应的正交规范特征向量。现欲求 A 的 p 个占优特征对 (λ_i, x_i) ($i = 1, 2, \dots, p$)。

设 $Q_0 = (q_1^{(0)}, q_2^{(0)}, \dots, q_p^{(0)})$ ，而 $Q_0^T Q_0 = I_p$ 。 $q_1^{(0)}, q_2^{(0)}, \dots, q_p^{(0)}$ 可看作 p 个占优特征向量的初始近似，也可以看作 p 维子空间的正交基。按幂法的思想，应以矩阵 A 同时左乘这些向量，得

$$V_1 = A Q_0, \quad (4.4.1)$$

但是就一般而言，矩阵 V_0 不再是列正交了，所以应对 V_0 的列正交化，这可以利用修正的 Gram-Schmidt 正交化算法来实现。于是得

$$V_0 = Q_1 R_1 \quad (4.4.2)$$

或

$$A Q_0 = Q_1 R_1. \quad (4.4.3)$$

这里 $Q_1 \in \mathbb{R}^{n \times p}$ 为列正交阵, $R \in \mathbb{R}^{p \times p}$ 为上三角阵。进而又可取 Q_1 作为新的迭代矩阵, 因为 Q_1 的列实际上是 V_0 列空间的正交基 (注意 R_1 为非异)。一般的迭代格式为:

$$\begin{cases} V_k = A Q_k, \\ V_k = Q_{k+1} R_{k+1} \end{cases} \quad (k=0, 1, \dots). \quad (4.4.4)$$

以上迭代格式常称为简单子空间迭代法。容易验证:

$$A^k Q_0 = Q_k R_k R_{k-1} \cdots R_1. \quad (4.4.5)$$

可以证明: 当 $k \rightarrow \infty$ 时, Q_k 的列将趋于 A 的 p 个占优正交规范特征向量, 但收敛速度很慢。

Q_k 的列收敛速度缓慢的原因在于: 在子空间 $\text{Span}(Q_k)$ 中, Q_k 的列未必是 A 的 p 个占优特征向量在空间 $\text{Span}(Q_k)$ 中的最优近似。于是, 我们必须研究以下问题: 设给定 $Q \in \mathbb{R}^{n \times p}$ 且 $Q^T Q = I_p$, 据此构成 p 维子空间 $\mathfrak{S}^p \equiv \text{Span}(Q)$ 。那末, 在该子空间范围中, A 的 p 个占优特征对的最优近似是什么?

如果 $p=1$, 即给定向量 q 且 $\|q\|_2 = 1$, 它作为 A 的近似特征向量, 对应的 A 的近似特征值 γ 应取何值为最优呢? 很自然, 如果 γ 值使残量 $\|Aq - q\gamma\|_2$ 为极小, 那末该 γ 值即为最优值。据广义的勾股定理知最优的 γ 使

$$r(q) = Aq - q\gamma \perp q$$

亦即

$$q^T (Aq - q\gamma) = 0$$

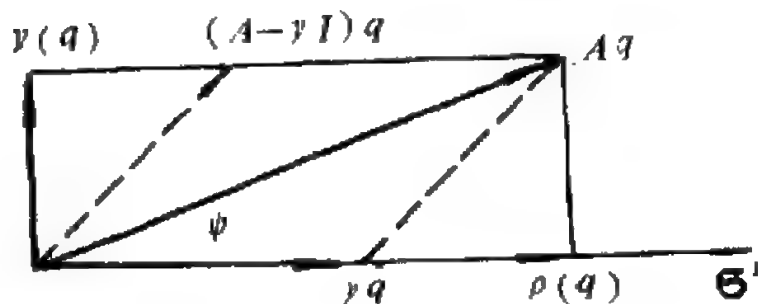


图 4.2 广义勾股定理

由此即可推知最优解即为 A 的 Rayleigh 商

$$r = r q^T q = q^T A q = \rho(q). \quad (4.4.6)$$

如果 $p > 1$, 则有以下定理

定理 4.4.1 设 $A \in R^{n \times n}$ 为对称阵, 给定 $Q \in R^{n \times p}$ 且 $Q^T Q = I_p$, 则对任何 $B \in R^{p \times p}$ 成立不等式

$$\|AQ - QB\|_p \leq \|AQ - QB\|_p \quad (p = 2, F), \quad (4.4.7)$$

其中 $H = Q^T A Q = \rho(Q)$.

证 因为

$$(AQ - QB)^T (AQ - QB) = Q^T A^2 Q - H^2 + (H - B)^T (H - B),$$

而

$$\|AQ - QB\|_F^2 = \text{tr}[(AQ - QB)^T (AQ - QB)],$$

所以

$$\|AQ - QB\|_F^2 \geq \|AQ - QH\|_F^2.$$

另一方面, 因为 $(H - B)^T (H - B)$ 为非负定, 故

$$\begin{aligned} \|AQ - QB\|_2^2 &= \max[\lambda((AQ - QB)^T (AQ - QB))] \\ &\geq \max[\lambda((AQ - QH)^T (AQ - QH))] \\ &= \|AQ - QH\|_2^2. \end{aligned}$$

以上定理说明, 虽然 $\mathcal{S}^p \equiv \text{span}(a)$ 不是 A 的不变子空间, 可是当取 B 为 A 的 Rayleigh 商矩阵 $\rho(Q)$ 时, 能使残量 $\|AQ - QB\|_p$ ($p = 2, F$) 极小.

现设 $\lambda(H): \theta_1 \leq \theta_2 \leq \dots \leq \theta_p$, 那末它们是不是由子空间 $\mathcal{S}^p \equiv \text{span}(Q)$ 推得的 A 的 p 个占优特征值的最优近似呢?

如果 \mathcal{S}^p 为 A 的 p 个占优特征向量张成的不变子空间, 则显然有

$$\lambda_{n-p+j} = \min_{\mathcal{S}' \subset \mathcal{S}^p} \max_{\substack{g \in \mathcal{S}' \\ g^T g = 1}} \frac{g^T A g}{g^T g} \quad (j = 1, \dots, p),$$

其中 \mathcal{S}' 为 \mathcal{S}^p 中任意 j 维子空间. 如果 \mathcal{S}^p 不是 A 的 p 个占优特征

向量张成的不变子空间，则很自然地定义：

$$\beta_j = \min_{\mathbf{G} \subset \mathbf{E}^p} \max_{\substack{\mathbf{g} \in \mathbf{G}^j \\ \mathbf{g} \neq 0}} \frac{\mathbf{g}^T \mathbf{A} \mathbf{g}}{\mathbf{g}^T \mathbf{g}} \quad (j = 1, \dots, p). \quad (4.4.8)$$

把它作为由 \mathbf{G}^j 推得的 λ_{n-p+j} 的最优近似是合理的。那末， β_j 与 θ_j 之间有什么关系呢？以下定理回答了这个问题。在证明定理之前，首先应注意子空间之间存在下列关系：

$$\begin{aligned} & \mathbf{R}^n \\ & \cup \\ & \mathbf{G}^p \equiv \mathbf{Q} \mathbf{R}^p \quad (\mathbf{Q} \in \mathbf{R}^{n \times p}) \\ & \cup \quad \cup \\ & \mathbf{G}^j \equiv \mathbf{Q} \mathcal{J}^j \equiv \mathbf{Q} \mathbf{G}^j \quad (\mathbf{G} \in \mathbf{R}^{p \times j}). \end{aligned}$$

请注意， \mathbf{G}^j 为 \mathbf{G}^p 中的 j 维子空间，也是 \mathbf{R}^n 中的 j 维子空间，而 \mathcal{J}^j 为 \mathbf{R}^p 中的 j 维子空间。由以上关系可见， $\mathbf{G}^j \subset \mathbf{R}^n$ 与 $\mathcal{J}^j \subset \mathbf{R}^p$ 之间存在一一对应的关系。

定理 4.4.2 $\beta_j = \lambda_j(\mathbf{H}) \equiv \theta_j \quad (j = 1, \dots, p)$

$$\begin{aligned} \text{证} \quad \beta_j &= \min_{\mathbf{G} \subset \mathbf{E}^p} \max_{\substack{\mathbf{g} \in \mathbf{G}^j \\ \mathbf{g} \neq 0}} \frac{\mathbf{g}^T \mathbf{A} \mathbf{g}}{\mathbf{g}^T \mathbf{g}} \quad (\mathbf{g} = \mathbf{Q} \mathbf{s}, \mathbf{s} \in \mathcal{J}^j, \\ & \quad \mathbf{g}^T \mathbf{g} = \mathbf{s}^T \mathbf{Q}^T \mathbf{Q} \mathbf{s} = \mathbf{s}^T \mathbf{s} \neq 0) \end{aligned}$$

$$= \min_{\mathcal{J} \subset \mathbf{R}^p} \max_{\substack{\mathbf{s} \in \mathcal{J}^j \\ \mathbf{s} \neq 0}} \frac{\mathbf{s}^T \mathbf{Q}^T \mathbf{A} \mathbf{Q} \mathbf{s}}{\mathbf{s}^T \mathbf{s}}$$

$$= \min_{\mathcal{J} \subset \mathbf{R}^p} \max_{\substack{\mathbf{s} \in \mathcal{J}^j \\ \mathbf{s} \neq 0}} \frac{\mathbf{s}^T \mathbf{H} \mathbf{s}}{\mathbf{s}^T \mathbf{s}}$$

$$= \theta_j.$$

由以上定理可见： \mathbf{H} 的特征值确是由子空间 \mathbf{G}^j 推得的 \mathbf{A} 的 p 个占优特征值的最优近似。

设 $\mathbf{H} \mathbf{g}_j = \theta_j \mathbf{g}_j \quad (j = 1, \dots, p)$ ， $\|\mathbf{g}_j\|_2 = 1$ 且 $\mathbf{g}_i^T \mathbf{g}_j = 0 \quad (i \neq j)$ ，

则

$$\begin{pmatrix} \theta_1 & & 0 \\ & \ddots & \\ 0 & & \theta_p \end{pmatrix} = G^T H G = (Q G)^T A (Q G), \quad (4.1.9)$$

这里 $G = (g_1, \dots, g_p)$. 因为 θ_i 可视为 λ_{n-p+i} 的最优近似, 所以取 $y_i = Q g_i$ 作为对应特征向量的最优近似也是合理的, 进一步的论述可参见本章习题 4.8. θ_i 和 y_i 分别称为矩阵 A 在子空间 $\text{Span}(Q)$ 上的 Ritz 值和 Ritz 向量. Ritz 值和 Ritz 向量总称为矩阵 A 在子空间 $\text{Span}(Q)$ 上的 Rayleigh-Ritz 近似.

综上所述即可给出 Rayleigh-Ritz 过程, 以下简称 RR 过程.

RR 过程给定对称阵 $A \in \mathbb{R}^{n \times n}$ 及 p 维子空间 $\mathcal{S} \equiv \{Sx \mid x \in \mathbb{R}^p\}$, 其中 $S \in \mathbb{R}^{n \times p}$ 而 $\text{rank}(S) = p$. 利用 \mathcal{S} 欲求 A 的 p 个占优特征对的最优近似 (θ_i, y_i) ($i = 1, \dots, p$).

(i) 调用修正的 Gram-Schmidt 子程序, 正交规范化 S 的列即

$$S = QR;$$

(ii) 调用子程序 OP 构成 AQ ;

(iii) 构成 Q 的 Rayleigh 商

$$H = \rho(Q) = Q^T (AQ);$$

(iv) 调用 Jacobi 子程序, 作谱分解

$$H g_i = g_i \theta_i, \quad (i = 1, \dots, p);$$

(v) 形成 Ritz 向量

$$y_i = Q g_i, \quad (i = 1, \dots, p).$$

上述过程中 OP 子程序的功能是: 矩阵 A 毋需用一个 $n \times n$ 数组给出, 而给出了对任意向量 v 计算 $A*v$ 的规则. 在该子程序中, 矩阵 A 的存贮以 $A*v$ 的计算都充分利用 A 的稀疏性.

至此，我们即可将简单子空间迭代法和 Rayleigh-Ritz 过程结合起来，产生新的子空间迭代算法。

算法 4.4.1 设给定对称阵 $A \in R^{n \times n}$ 及 $S_0 \in R^{n \times l}$ ，而 $\text{rank}(S_0) = p$ ，从子空间 $\mathfrak{S} = \text{span}(S_0)$ 出发，求 A 的 p 个占优特征对。

1) 对 $k = 1, 2, \dots$,

1.1) 调用 OP 子程序形成

$$C_k = AS_{k-1}.$$

1.2) 调用修正的 Gram-Schmidt 子程序，正交规范化 C_k 的列，即

$$C_k = Q_k R_k,$$

1.3) 调用 OP 子程序构成 AQ_k ，

1.4) 构成 Rayleigh 商矩阵

$$\hat{H}_k = Q_k^T A Q_k,$$

1.5) 调用 Jacobi 子程序，谱分解

$$\hat{H}_k = G_k \Theta G_k^T, \quad (\Theta = \text{diag}(\theta_1, \dots, \theta_p)),$$

1.6) 形成 Ritz 向量

$$S_k = Q_k G_k,$$

1.7) 检验是否达到收敛标准，若已达到收敛标准即转3)。

2) NEXT k 。

3) 输出 (θ_i, y_i) ($i = 1, \dots, p$)。

4) END。

一般而言，算法 4.4.1 中 G_k 为一系列正交矩阵的乘积，
设

$$G_k = p_1 \cdots p_l$$

那末语句1.5) 和语句1.6) 可同时进行如下：

1) 置 $S_k = Q_k$;

2) 对 $i = 1, \dots, l$, 置 $S_k = S_k P_i$.

这样矩阵 G_k 就毋需明显形成, 在整个计算过程中只要一个 $n \times p$ 数组就是以应付 S 、 C 以及 Q 的计算了. 其次, 读者应注意, 在算法 4.4.1 中, S_k 的列是子空间 $A^k \mathcal{S}$ 的基, C_k 的列也是子空间 $A^k \mathcal{S}$ 的基. 但前者应较后者收敛于 A 的 p 个占优特征向量要快. 当然迭代一次的计算量相对而言较高. 特别是当矩阵阶数较高时, 在执行语句 1.3) 时要额外调用 **OP** 子程序 p 次, 显然算法 4.4.1 的计算量要大大超过简单子空间迭代法的计算量. 为此, 以下我们来讨论如何避免额外调用 p 次 **OP** 子程序.

因为考虑到 A^{-2} 与 A 具有相同的特征向量, 所以也可针对 A^{-2} 来实施 Rayleigh-Ritz 过程.

如同算法 4.4.1, 首先计算 $C_k = AS_{k-1}$, 然后设法形成新基 $C_k F_k$, 这里 $F_k \in \mathbb{R}^{p \times p}$ 应满足以下二个条件:

$$\textcircled{1} (C_k F_k)^T (C_k F_k) = I_p, \quad (\text{正交性}).$$

$$\textcircled{2} (C_k F_k)^T A^{-2} (C_k F_k) = D_k^{-2}, \quad (\text{Ritz 向量性质}).$$

其中 $D_k = \text{diag}(\delta_1, \dots, \delta_p)$ 为 p 阶对角阵. ②式亦可写成

$$(F_k^T S_{k-1}^T A) A^{-2} (AS_{k-1} F_k) = F_k^T F_k = D_k^{-2}. \quad (4.4.10)$$

由此可见, $F_k D_k$ 为正交阵. 这样据①式及 (4.4.10) 即有

$$\begin{aligned} C_k^T C_k &= F_k^{-T} F_k^{-1} = (F_k^{-T} D_k^{-1}) D_k^2 (D_k^{-1} F_k^{-1}) \\ &= (F_k D_k) D_k^2 (F_k D_k)^T. \end{aligned} \quad (4.4.11)$$

由此可见, F_k 和 D_k 即可由 $C_k^T C_k$ 的谱分解确定.

算法 4.4.2 设给定对称阵 $A \in \mathbb{R}^{n \times n}$ 及 $S_0 \in \mathbb{R}^{n \times p}$ 而 $\text{rank}(S_0) = p$. 从子空间 $\mathcal{S} = \text{span}(S_0)$ 出发, 欲求 A 的 p 个占优特征对.

1) 对 $k = 1, 2, \dots$,

1.1) 调用 **OP** 子程序形成

$$C_k = A s_{i-1},$$

1.2) 构成矩阵

$$\dot{H}_k = C_k^T C_k,$$

1.3) 调用Jacobi子程序, 作谱分解

$$\dot{H}_k = B_k D_k^2 B_k^T,$$

1.4) 形成

$$S_k = C_k B_k D_k^{-1},$$

1.5) 检验是否达到收敛标准, 若已达到收敛标准即转3).

2) NEXT k .

3) 输出 (δ_i, y_i) ($i = 1, \dots, p$), 这里 y_i 是 S_k 的第 i 列.

4) END.

应注意算法4.4.2中的 S_k 和算法4.4.1中的 S_k 是不同的, 而且算法4.4.2中的 S_k 并不能表成一系列正交阵的乘积. 为此对算法4.4.2仍须作进一步修改. 考虑到

$$\dot{H}_k = C_k^T C_k = R_k^T Q_k^T Q_k R_k = R_k^T R_k = L_k R_k. \quad (4.4.12)$$

如果记 $H_k \equiv R_k L_k$, 那末据LR算法知, H_k 较 \dot{H}_k 更接近对角阵. 但据(4.4.12)应有

$$\begin{aligned} H_k &\equiv R_k L_k \equiv R_k R_k^T = R_k \dot{H}_k R_k^{-1} \\ &= R_k B_k D_k^2 B_k^T R_k^{-1} \\ &= (R_k B_k D_k^{-1}) D_k^2 (D_k B_k^T R_k^{-1}). \end{aligned} \quad (4.4.13)$$

由上式知, H_k 与 \dot{H}_k 相似, 所以用 H_k 替代 \dot{H}_k 是合理的. 其次, 不难验证, $P_k \equiv R_k B_k D_k^{-1}$ 是 H_k 的正交特征矩阵, 于是算法4.4.2中的 S_k 即可表为

$$S_k \equiv C_k B_k D_k^{-1} = Q_k R_k B_k D_k^{-1} = Q_k P_k.$$

这样 S_k 就有可能表成一系列正交阵的乘积了. 综上所述, 即可构成如下 RITZIT 算法.

算法4.4.3 设给定对称阵 $A \in R^{n \times n}$ 及 $S_0 \in R^{n \times p}$, 而

$\text{rank}(S_0) = p$, 从子空间 $\mathcal{S} = \text{span}(S_0)$ 出发, 求 A 的 p 个占优特征对,

1) 对 $k = 1, 2, \dots$,

1.1) 调用 OP 子程序形成

$$C_k = AS_{k-1},$$

1.2) 调用修正的 Gram-Schmidt 子程序, 正交规范化 C_k 的列, 即

$$C_k = Q_k R_k,$$

1.3) 形成矩阵

$$H_k = R_k R_k^T,$$

1.4) 调用 Jacobi 子程序, 作谱分解

$$H_k = P_k D_k^A P_k^T,$$

1.5) 形成 Ritz 向量

$$S_k = Q_k P_k,$$

1.6) 检验是否达到收敛标准, 若已达到收敛标准, 即转3).

2) NEXT k .

3) 输出 (δ_i, y_i) ($i = 1, \dots, p$).

4) END.

在一定的假设下, 可以证明子空间迭代法的收敛性。我们不加证明引入以下定理。

定理4.4.3 设 $A \in R^{n \times n}$ 为对称正定阵, 其特征值及其对应的特征向量分别为

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_{n-p} < \lambda_{n-p+1} \leq \dots \leq \lambda_n$$

和

$$s_1, s_2, \dots, s_n.$$

设 $S \in R^{n \times p}$ 且 $S^T S = I_p$, 若 $Z^T S$ 为非异, 则由算法4.4.3产生的 Ritz 向量当 $k \rightarrow \infty$ 时有

$$\sin(\mathbf{y}_i^{(k)}, \mathbf{z}_{n-p+i}) = O\left[\left(\frac{\lambda_{n-p}}{\lambda_{n-p+i}}\right)^k\right], \quad (4.4.14)$$

$$(i = 1, \dots, p)$$

这里 $\mathbf{Z} = (\mathbf{z}_n, \mathbf{z}_{n-1}, \dots, \mathbf{z}_p)$ 。

据以上定理知，当 $i = p$ 时有

$$\sin(\mathbf{y}_p^{(k)}, \mathbf{z}_n) = O\left[\left(\frac{\lambda_{n-p}}{\lambda_n}\right)^k\right]. \quad (4.4.15)$$

而一般而言， $\left|\frac{\lambda_{n-p}}{\lambda_n}\right| < \left|\frac{\lambda_{n-1}}{\lambda_n}\right|$ ，所以利用 RR 过程等措施改进后的子空间迭代法，显然要比简单子空间迭代法有效得多。

由 (4.4.14) 还可以看出：在相同的迭代次数前提下，愈是排在后面的特征值及其对应的特征向量，它们的计算精度愈是高。因此，在实际计算时，往往把 p 取得比实际所求的个数 r 稍大一点。有的文献推荐，取 $p = \min(2r, r + 8)$ ，这额外的 $p - r$ 个向量称为护卫向量。显然，对于固定的 r ， p 取得愈大，迭代次数就愈小。但这样每次迭代所花费的工作量就较大。因此这里有个权衡得失的问题，这需在实际计算中去探索。至于 \mathbf{S}_0 的选取，可以先取 p 个由随机数构成的向量，然后将它们正交规范化。

最后我们指出，如果要求 \mathbf{A} 前面 p 个按模小的特征值，只要在以上算法中将 \mathbf{A} 换成 \mathbf{A}^{-1} 即可。

下面我们来讨论对称广义特征值问题的子空间迭代法。正如 §1 指出那样，对称广义特征值问题

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{B}\mathbf{x} \quad (4.4.16)$$

可化为标准特征值问题

$$\begin{aligned} \tilde{\mathbf{A}}\mathbf{y} &= \lambda\mathbf{y}, \\ \tilde{\mathbf{A}} &= \mathbf{L}^{-1}\mathbf{A}\mathbf{L}^{-T}, \\ \mathbf{y} &= \mathbf{L}^T\mathbf{x}, \\ \mathbf{B} &= \mathbf{L}\mathbf{L}^T. \end{aligned} \quad (4.4.17)$$

于是算法 4.4.3 和 (4.4.17) 相结合, 即可得到对称广义特征值问题的子空间迭代算法.

现在我们来分析, 用子空间迭代法求解对称广义特征值问题的独到之处. 对于对称广义特征值问题而言, 算法 4.4.3 中语句 1.1) 实际上应执行

$$C_k = L^{-1} A L^{-T} S_{k-1}, \quad (4.4.18)$$

而

$$S_{k+1} = L^T x_{k+1}. \quad (4.4.19)$$

实际上, 以上运算可以分以下二步来完成:

(i) 调用 **OP** 子程序计算

$$Z_k = A x_{k+1}; \quad (4.4.20)$$

(ii) 用向后回代法求解

$$L C_k = Z_k. \quad (4.4.21)$$

如果 A, B 都是稀疏矩阵, 那末在执行以上几步运算时, 就可充分利用它们的稀疏性质. 首先, 在执行迭代之前, 对 B 作 **Cholesky** 分解, 这可利用后面第六章所介绍的对称变带宽矩阵 **Cholesky** 分解方法. 求得 L 后, 在迭代过程中只要反复求解 (4.4.21) 形状的方程组集, 每次迭代时仅 Z_k 在变化. 其次, 对矩阵 A 亦可利用第六章介绍的任何稀疏矩阵压缩存贮技巧, 以提高 **OP** 子程序的效率.

算法 4.4.4 设给定对称广义特征值问题 $Ax = \lambda Bx$, 设 $S_0 \in R^{n \times p}$ 且 $S_0^T S_0 = I_p$, 求问题的前 p 个占优特征值及其对应的特征向量.

1) 对 B 作 **Cholesky** 分解

$$B = LL^T.$$

2) A 按某种对称稀疏阵存贮方法存贮

3) 求解

$$L^T X_0 = S_0.$$

4) 对 $k = 1, 2, \dots$,

4.1) 调用 OP 子程序计算

$$Z_k = AX_{k-1},$$

4.2) 求解

$$LC_k = Z_k,$$

4.3) 调用修正的 Gram-schmidt 子程序, 作

$$C_k = Q_k R_k,$$

4.4) 形成

$$H_k = R_k R_k^T,$$

4.5) 调用 Jacobi 子程序, 作谱分解

$$H_k = P_k D_k^2 P_k^T,$$

4.6) 形成

$$S_k = Q_k P_k,$$

4.7) 求解

$$L^T X_k = S_k,$$

4.8) 检验是否达到收敛标准, 若已达到收敛标准即转

6).

5) NEXT k .

6) 输出 (δ_i, x_i) ($i = 1, \dots, p$).

7) END

第四章 习 题

4.1 求解以下广义特征值问题 $Ax = \lambda Bx$,

$$(i) \quad A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix},$$

$$(ii) \quad A = \begin{pmatrix} 2 & 2 \\ 2 & 4 \end{pmatrix}, \quad B = \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix},$$

$$(iii) A = \begin{pmatrix} 5 & -2 & 0 \\ -2 & 3 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix}.$$

4.2 将以下对称广义特征值问题 $Ax = \lambda Bx$ 化为标准特征值问题 (4.1.6):

$$(i) A = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}, \quad B = \begin{pmatrix} 2 & 1 \\ 1 & 4 \end{pmatrix},$$

$$(ii) A = \begin{pmatrix} -2 & 1 & 0 & 0 \\ 1 & -2 & 1 & 0 \\ 0 & 1 & -2 & 1 \\ 0 & 0 & 1 & -2 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

4.3 设 $A, B \in R^{n \times n}$ 且 B 为非异, 试证存在酉阵 U 和 V , 使 A, B 同时酉等价于上三角阵 \tilde{A}, \tilde{B} :

$$\tilde{A} = UHAU, \quad \tilde{B} = UHBV.$$

4.4 设 $A, B \in R^{n \times n}$ 为对称矩阵且 B 为正定, $Ax = \lambda Bx$ 的特征值为 $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. 试证:

$$\lambda_j = \min_{\substack{\dim(\mathcal{J})=j \\ x \in \mathcal{J} \\ x \neq 0}} \max \frac{x^T A x}{x^T B x} \quad (j = 1, \dots, n).$$

这里 \mathcal{J} 为 R^n 中的 j 维子空间.

4.5 设 $A, B, C \in R^{n \times n}$ 皆为对称正定, 试证二次特征值问题:

$$(A - \lambda B - \lambda^2 C)x = 0.$$

(i) 等价于标准特征值问题:

$$(G - \lambda I_{2n})y = 0$$

其中

$$G = \begin{pmatrix} O & I_n \\ C^{-1}A & -C^{-1}B \end{pmatrix}, \quad y = \begin{pmatrix} x \\ \hat{x} \end{pmatrix}, \quad \hat{x} = \lambda x.$$

(ii) 等价于对称广义特征值问题

$$(E - \lambda F)y = 0$$

其中

$$E = \begin{pmatrix} C & O \\ O & A \end{pmatrix}, \quad F = \begin{pmatrix} O & C \\ C & B \end{pmatrix}, \quad y = \begin{pmatrix} \hat{x} \\ x \end{pmatrix}, \quad \hat{x} = \lambda x.$$

4.6 设对称特征值问题 $A - \lambda B$, 它对应的 Rayleigh 商为

$$\rho(x) = \frac{x^T A x}{x^T B x}, \quad x \neq 0.$$

试证:

$$(i) \rho(ax) = \rho(x), \quad a \neq 0.$$

$$(ii) \lambda_1 \leq \rho(x) \leq \lambda_n \quad (\forall \|x\|_2 = 1),$$

$$(iii) \nabla \rho(x) = 2(Ax - \rho(x)Bx) / x^T Bx,$$

$$(iv) \|(A - \sigma B)x\|_{A^{-1}}^2 \geq \|Ax\|_{A^{-1}}^2 - |\rho(x)|^2 \|Bx\|_{A^{-1}}^2,$$

等号成立当且仅当 $\sigma = \rho(x)$.

4.7 设 $p = F$, 试证 (4.4.7) 式等号成立当且仅当 $B = H$; 设 $p = 2$, 试构造矩阵 A , Q 以及二阶阵 B , 使 $B \neq H$ 但 (4.4.7) 式等号成立.

4.8 设 $A \in R^{n \times n}$ 为对称阵, \mathcal{J}^p 为 R^n 中的 p 维子空间, 试证对任何 $x \in \mathcal{J}^p$, $Ax - \rho(x)x \perp \mathcal{J}^p$ 当且仅当 x 为矩阵 A 在子空间 \mathcal{J}^p 上的 Ritz 向量.

4.9 利用算法 4.3.1 确定以下对称矩阵:

$$(i) \begin{pmatrix} 1 & 0 & 1 & 2 \\ 0 & -1 & -2 & 0 \\ 1 & -2 & -2 & 3 \\ 2 & 0 & 3 & 6 \end{pmatrix}, \quad (ii) \begin{pmatrix} 1 & 2 & 1 & 2 \\ 2 & -1 & -2 & 0 \\ 1 & -2 & -2 & 3 \\ 2 & 0 & 3 & 6 \end{pmatrix}$$

各阶首主子式的符号.

4.10 编制算法 4.4.3 的电算程序, 利用该程序计算矩阵

$$A = \begin{pmatrix} -1 & 0 & 3 & 0 \\ 0 & -1 & 0 & 0 \\ 3 & 0 & -1 & 0 \\ 0 & 0 & 0 & -3 \end{pmatrix}$$

的前二个占优特征对, 精度自定.

4.11 编制算法 4.4.4 的电算程序, 利用该程序计算 $Ax = \lambda Bx$ 的前二个占优特征对, 精度自定。矩阵 A, B 为

$$A = \begin{pmatrix} 10 & 2 & 3 & 1 & 11 \\ 2 & 12 & 1 & 2 & 1 \\ 3 & 1 & 11 & 1 & -1 \\ 1 & 2 & 1 & 9 & 1 \\ 1 & 1 & -1 & 1 & 15 \end{pmatrix}, \quad B = \begin{pmatrix} 12 & 1 & -1 & 2 & 1 \\ 1 & 14 & 1 & -1 & 1 \\ -1 & 1 & 16 & -1 & 1 \\ 2 & -1 & -1 & 12 & -1 \\ 1 & 1 & 1 & -1 & 11 \end{pmatrix}.$$

参 考 文 献

1. Stewart, G. W., (1973), *Introduction to Matrix Computation*, New York, Academic Press.
2. Jennings, A., (1977), *Matrix Computation for Engineers and Scientists*, New York, John Wiley.
3. Wilkinson, J. H., (1965) *The Algebraic Eigenvalue Problem*, New York: Oxford Univ. Press.
4. Wilkinson, J. H. & Reinsch, C. H. (1971), *Handbook for Automatic Computation. Linear Algebra*, Vol. 2, New York: Springer-Verlag.
5. Gourlay, A. R. & Watson, G. A., (1973), *Computational Methods for Matrix Eigenproblems*, New York: John Wiley.
6. Parlett, B. N., (1980), *The Symmetric Eigenvalue Problem*, Englewood Cliffs, N. J.: Prentice-Hall.
7. 曹志浩等矩阵计算和方程求根。
8. 李大潜等有限元素法续讲。
9. Moler, C. B. & Stewart, G. W. (1973), "An algorithm for generalized matrix eigenvalue Problems" SIAM. J. Numer. Anal. 10. 241-256.
10. Stewart, G. W. (1976), "Simultaneous iteration for computing

invariant subspaces of non-Hermitian matrices" Numer. Math. 25, 123-136.

11. Clint, M., & Jennings, A. (1971) "*A simultaneous iteration method for the unsymmetric eigenvalue Problems*" J. Inst. Maths, Applies, 8, 111-121.

第五章 线性最小二乘法

在测量平差、数据逼近以及正态分布有关的统计等问题中，都需要考虑确定一个向量 \mathbf{x} ，使函数

$$\rho^2(\mathbf{x}) = \|\mathbf{b} - A\mathbf{x}\|^2$$

达到最小的问题。其中 A 是一个 $m \times n$ 的矩阵， \mathbf{b} 是 m 维列向量，这个问题称为**最小二乘法问题**。本章将通过对最小二乘法的讨论来介绍 $m \times n$ 矩阵的奇异值分解以及广义逆矩阵的某些概念和有关的递推算法。

§1 线性最小二乘法问题

在某些实际问题中碰到的线性方程组 $A\mathbf{x} = \mathbf{b}$ 往往无解，即找不到一个向量 \mathbf{x}_0 ，使 $A\mathbf{x}_0 = \mathbf{b}$ ，此时只能将问题合理地改为：寻找列向量 \mathbf{x} 使 $\rho^2(\mathbf{x}) = \|\mathbf{b} - A\mathbf{x}\|^2$ 达到最小。即使剩余向量的模的平方

$$\|\mathbf{r}\|^2 = \|\mathbf{b} - A\mathbf{x}\|^2 \quad (5.1.1)$$

达到最小。由于剩余向量 \mathbf{r} 是线性地依赖于 \mathbf{x} 的，所以我们称这个问题为线性最小二乘法问题。而把使(5.1.1)右端达到最小的 \mathbf{x} 叫做

$$A\mathbf{x} = \mathbf{b} \quad (5.1.2)$$

的最小二乘解。

现在来研究最小二乘解的存在性以及它的唯一性条件。为了以后的需要，我们首先对长方形矩阵的一些基本性质做如下讨论。

设长方形矩阵 $A \in R^{m \times n}$, $\|A\|_2$ 的定义可以仿方阵的 2 范数定义, 定义为

$$\|A\| = \|A\|_2 = \text{Max}_{x \neq 0} (\|Ax\|_2 / \|x\|_2), \quad (5.1.3)$$

可以证明它具有下列性质:

- 1) $\|A\|^2 = A^T A$ 的最大特征值;
- 2) $\|A^T\| = \|A\|$;
- 3) $\|A^T A\| = \|A\|^2$.

对于满足方程组

$$Ax = 0$$

的所有 x 组成的空间称为 A 的零空间, 记为 $N(A)$, 它的维数记为 $\text{null}(A)$. 以后也用 $R(A)$ 来表示 A 的列向量张成的线性空间.

定理 5.1.1 使 (5.1.1) 达到最小的线性最小二乘解总是存在的, 而且解唯一的充要条件是 $\text{null}(A) = 0$.

证明 根据线性代数的知识, b 可唯一地分解为

$$b = b_1 + b_2,$$

其中 $b_1 \in R(A)$, $b_2 \in R^\perp(A)$, $R^\perp(A)$ 是 A 的正交补空间. 因为 $Ax \in R(A)$, 所以 $-Ax \in R(A)$, 故 $b - Ax$ 有唯一分解式

$$b - Ax = b_1 + b_2 - Ax, \quad (5.1.4)$$

其中 $b_1 - Ax \in R(A)$, $b_2 \in R^\perp(A)$. 从而可以得到

$$\|b - Ax\|^2 = \|b_1 - Ax\|^2 + \|b_2\|^2. \quad (5.1.5)$$

由于 $b_1 \in R(A)$, 所以线性方程组

$$Ax = b_1 \quad (5.1.6)$$

必有解, 设解为 \bar{x} , 即 $b_1 - A\bar{x} = 0$, 于是对任意向量 x 均有

$$\|b - Ax\|^2 = \|b_1 - Ax\|^2 + \|b_2\|^2 \geq \|b_2\|^2 = \|b - A\bar{x}\|^2.$$

这就证明 \bar{x} 是 (5.1.2) 的线性最小二乘解, 从而证明了线性最小二乘法问题总有解.

设 \bar{x}_0 是方程组 (5.1.6) 的任一解, 则 (5.1.6) 的所有解的集合可以表示为 $\bar{x}_0 + N(A)$. 由此可得, 如果 $\text{null}(A) = 0$, 则 (5.1.6) 的解唯一, 否则, 又可得 (5.1.6) 的解不唯一, 从而证明了 (5.1.6) 的解唯一存在的充要条件是 $\text{null}(A) = 0$.

定理 5.1.2 (5.1.2) 的线性最小二乘解是方程组

$$A^T A \bar{x} = A^T \bar{b} \quad (5.1.7)$$

的解, 反之, 方程组 (5.1.7) 的解是 (5.1.2) 的最小二乘解.

证明 设 $\bar{b} - A\bar{x} = \bar{b}_1 + \bar{b}_2 - A\bar{x}$, 其中 $\bar{b}_1 \in R(A)$, $\bar{b}_2 \in R^\perp(A)$, 若 \bar{x} 是 (5.1.2) 的最小二乘解, 则必有 $A\bar{x} = \bar{b}_1$, 因此, $\bar{b} - A\bar{x} = \bar{b} - \bar{b}_1 = \bar{b}_2 \in R^\perp(A)$. 又由于 $R^\perp(A) = N(A^T)$, 故 $\bar{b} - A\bar{x} = \bar{b}_2 \in N(A^T)$, 即 $A^T(\bar{b} - A\bar{x}) = 0$, 从而证明了 \bar{x} 是方程组 (5.1.7) 的解.

反之, 若 \bar{x} 是方程组 (5.1.7) 的一个解, 则必有

$$A^T \bar{b} - A^T A \bar{x} = A^T(\bar{b} - A\bar{x}) = 0.$$

因此,

$$\bar{b} - A\bar{x} \in R^\perp(A).$$

但是, $\bar{b} - A\bar{x} = \bar{b}_2 + \bar{b}_1 - A\bar{x}$, 其中 $\bar{b}_1 - A\bar{x} \in R(A)$, $\bar{b}_2 \in R^\perp(A)$, 根据 $\bar{b} - A\bar{x}$ 在 $R(A)$ 、 $R^\perp(A)$ 上分解的唯一性得到 $\bar{b}_1 - A\bar{x} = 0$, 于是 \bar{x} 是 (5.1.2) 的线性最小二乘解.

定义 5.1.1 方程组 $A^T A \bar{x} = A^T \bar{b}$ 称为线性最小二乘法问题 (5.1.1) 的法方程组.

由定理 5.1.1 和定理 5.1.2 可知, 当 $m \times n$ 矩阵 A 的零空间 $N(A)$ 的维数满足 $\text{null}(A) = 0$ 时, 线性最小二乘解才是唯一的, 即 A 的列线性无关时才有唯一的 \bar{x} 使 (5.1.1) 达到最小. 这时必有 $m \geq n$, 而且 $A^T A$ 是非奇异阵, 因此, $A^T A$ 是对称正定

的。所以，在 $\text{null}(A) = 0$ 时，求线性最小二乘解问题可以化为解系数阵为对称正定阵的法方程组(5.1.7)的问题。至于 $\text{null}(A) \neq 0$ 的情况，线性最小二乘法问题的求解就变得复杂了。我们将分别进行讨论。

§2 法方程组的解法

本节主要讨论方程组(5.1.2)的系数阵 A 的列向量线性无关时(简称列满秩阵)的最小二乘解。此时， $A^T A$ 是非奇异阵，从而方程组(5.1.7)的解可表为

$$x = (A^T A)^{-1} A^T b.$$

若令

$$\bar{A}^+ = (A^T A)^{-1} A^T, \quad (5.2.1)$$

则称矩阵 \bar{A}^+ 为 A 的Moore-Penrose广义逆矩阵。因此，线性最小二乘解可以表为

$$x = \bar{A}^+ b. \quad (5.2.2)$$

2.1 法方程组的性态

为了选择法方程组的合适的求解方法，应对法方程组(5.1.7)的特点进行研究。

形成法方程组的系数阵，需做矩阵乘积 $A^T A$ ，这不仅要花费很大的工作量，而且由于舍入误差的影响还有可能破坏 $A^T A$ 的正定性。其次，线性方程组(5.1.7)往往是病态的。例如，给定12个点

$(-6, y_1), (-5, y_2), \dots, (-1, y_6), (1, y_7), \dots, (5, y_{11}), (6, y_{12})$ 要求确定一个8次多项式

$$y(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_8 x^8$$

使 $\sum (y(x_i) - y_i)^2$ 达到最小，这需要解方程组

$$C\mathbf{a} = \mathbf{y},$$

其中

$$C = \begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^8 \\ 1 & x_2 & x_2^2 & \cdots & x_2^8 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & x_{12} & x_{12}^2 & \cdots & x_{12}^8 \end{pmatrix},$$

而 $\alpha = (\alpha_0, \alpha_1, \cdots, \alpha_8)^T$, $y = (y_1, y_2, \cdots, y_{12})^T$, 相应的法方程组为

$$C^T C \alpha = C^T y. \quad (5.2.3)$$

设 $C^T C = (c_{ij})_{9 \times 9}$, 则方程组 (5.2.3) 的条件数估计为

$$\text{Cond}(C^T C) = \frac{C^T C \text{ 的最大特征值}}{C^T C \text{ 的最小特征值}} \geq 4.96 \times 10^{11}.$$

显然, 本例中的法方程组 (5.2.3) 是“病态”的。因而, 用一般的方法去解 (5.2.3) 将得不到满意的结果。

为什么法方程组 (5.1.7) 往往是病态的呢? 现在仅就方程组 (5.1.2) 的右端有扰动 δb 的情况来进行讨论。

设方程组 (5.1.2) 的右端由 b 变为 $b + \delta b$, 此时 (5.1.7) 的解变为 $x + \delta x$, 且

$$A^T A (x + \delta x) = A^T (b + \delta b),$$

于是有 $A^T A \delta x = A^T \delta b$, 即

$$\delta x = (A^T A)^{-1} A^T \delta b = \bar{A}^+ \delta b. \quad (5.2.4)$$

定理 5.2.1 若 (5.1.2) 的矩阵 A 是列满秩阵, 且 $b_1 \neq 0$, 则

$$\frac{\|\bar{A}^+ \delta b\|}{\|\bar{A}^+ b\|} \leq \text{Cond}(A) \frac{\|\delta b_1\|}{\|b_1\|}, \quad (5.2.5)$$

其中 $\text{Cond}(A) = \|\bar{A}^+\| \|A\|$.

证明 令 δb 在 $R(A)$ 和 $R^\perp(A)$ 上的分解为

$$\delta b = \delta b_1 + \delta b_2,$$

其中 $\delta \mathbf{b}_1 \in R(A)$, $\delta \mathbf{b}_2 \in R^\perp(A) = N(A^T)$, 则

$$\bar{A}^+ \delta \mathbf{b} = \bar{A}^+ (\delta \mathbf{b}_1 + \delta \mathbf{b}_2) = \bar{A}^+ \delta \mathbf{b}_1 + (A^T A)^{-1} A^T \delta \mathbf{b}_2 = \bar{A}^+ \delta \mathbf{b}_1,$$

因此

$$\|\bar{A}^+ \delta \mathbf{b}\| = \|\bar{A}^+ \delta \mathbf{b}_1\| \leq \|\bar{A}^+\| \|\delta \mathbf{b}_1\|. \quad (5.2.6)$$

据(5.2.2)和(5.1.6)可得

$$A \bar{A}^+ \mathbf{b} = \mathbf{b}_1,$$

于是

$$\begin{aligned} \|\mathbf{b}_1\| &= \|A \bar{A}^+ \mathbf{b}\| \leq \|A\| \|\bar{A}^+ \mathbf{b}\| \\ \|\bar{A}^+ \mathbf{b}\| &\geq \|\mathbf{b}_1\| / \|A\|. \end{aligned} \quad (5.2.7)$$

由于 $\|\mathbf{b}_1\| \neq 0$, 用(5.2.7)除(5.2.6)得到

$$\frac{\|\bar{A}^+ \delta \mathbf{b}\|}{\|\bar{A}^+ \mathbf{b}\|} \leq \|\bar{A}^+\| \|A\| \frac{\|\delta \mathbf{b}_1\|}{\|\mathbf{b}_1\|} = \text{Cond}(A) \frac{\|\delta \mathbf{b}_1\|}{\|\mathbf{b}_1\|}.$$

(5.2.2), (5.2.4), (5.2.5)表明, 当(5.1.2)的右端有扰动 $\delta \mathbf{b}$ 时, 最小二乘解 \mathbf{x} 的相对误差界为

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \text{Cond}(A) \frac{\|\delta \mathbf{b}_1\|}{\|\mathbf{b}_1\|}.$$

该式的右端是 \mathbf{b}_1 的相对误差, 因而, 当(5.1.2)的右端 \mathbf{b} 变化时, 解的变化主要依赖于 \mathbf{b}_1 的变化.

定理5.2.2 若阵 $A \in R^{m \times n}$ 是列满秩阵, 则

$$\text{Cond}(A^T A) = (\text{Cond}(A))^2, \quad (5.2.8)$$

其中 $\text{Cond}(A^T A) = \|A^T A\| \|(A^T A)^{-1}\|$.

证明 因为 $\|A\|^2 = \|A^T A\| = \|A A^T\|$, 且

$$\bar{A}^+ = (A^T A)^{-1} A^T,$$

所以

$$\begin{aligned} \|\bar{A}^+\|^2 &= \|\bar{A}^+ (\bar{A}^+)^T\| \\ &= \|(A^T A)^{-1} A^T A (A^T A)^{-1}\| = \|(A^T A)^{-1}\| \end{aligned}$$

故

$$\begin{aligned}\text{Cond}(A^T A) &= \|A^T A\| \|(A^T A)^{-1}\| = \|A\|^2 \|\bar{A}^+\|^2 \\ &= (\text{Cond}(A))^2.\end{aligned}$$

定理5.2.2指出, 法方程组(5.1.7)的系数矩阵的条件数是原方程组(5.1.2)的系数矩阵 A 的条件数的平方, 即原来不太病态的问题, 化为法方程组后, 变得“惊人”地病态了. 因此必须采用具有更高数值稳定性的方法去求解.

2.2 正交化方法

线性代数中所介绍的 **Cholesky** 分解是把一个对称正定阵进行三角分解, 再将它应用于对 $A^T A$ 的分解上, 即得

$$A^T A = R^T R, \quad (5.2.9)$$

其中 R 是 n 阶非奇异的上三角阵, 且规定对角元取正号. 此时分解式(5.2.9)是唯一的. 如果令 $Q = AR^{-1}$, 则

$$A = QR, \quad (5.2.10)$$

而且 $Q^T Q = R^{-T} A^T A R^{-1} = R^{-T} R^T R R^{-1} = I$, 即 Q 为正交阵.

对于列满秩阵 A , 其QR分解可由Gram-Schmidt 正交化过程实现. 并把这个正交化过程简称为G-S过程. 设 $A = [a_1, a_2, \dots, a_n]$, $Q = [q_1, q_2, \dots, q_n]$, 则G-S过程的具体步骤如下:

1) 将 a_1 标准化, 作为 Q 的第1列 q_1 , 即

$$q_1 = a_1 / r_{11}, \quad r_{11} = \|a_1\|;$$

2) 对于 $k = 2, 3, \dots, n$, 按下列公式求出 Q 的相应列 q_k , 即

$$\begin{aligned}b_k &= a_k - \sum_{j=1}^{k-1} r_{jk} q_j, \quad r_{jk} = (q_j, a_k), \\ &\quad (j = 1, 2, \dots, k-1) \\ q_k &= b_k / r_{kk}, \quad r_{kk} = \|b_k\|.\end{aligned}$$

如此求出的 q_1, q_2, \dots, q_n 必为标准正交向量组, 这个过程可以写成如下形式:

$$\begin{aligned} \mathbf{a}_1 &= r_{11}\mathbf{q}_1, \\ \mathbf{a}_2 &= r_{12}\mathbf{q}_1 + r_{22}\mathbf{q}_2, \\ &\dots\dots\dots \\ \mathbf{a}_n &= r_{1n}\mathbf{q}_1 - r_{2n}\mathbf{q}_2 + \dots + r_{nn}\mathbf{q}_n. \end{aligned} \tag{5.2.11}$$

将(5.2.11)写成矩阵形式, 即为 $A = QR$, 其中

$$R = \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ 0 & r_{22} & \dots & r_{2n} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & r_{nn} \end{pmatrix}$$

Q 为正交阵.

将分解式(5.2.10)代入方程组(5.1.7), 得到

$$R^T Q^T Q R x = R^T Q^T b.$$

由于 R^T 是非奇异阵, 所以存在逆阵 $(R^T)^{-1}$, 用它左乘于上式两端, 得到

$$R\mathbf{x} = Q^T \mathbf{b}, \quad (5.2.12)$$

因此,求得 R 和 Q 后,用 $Q^T b$ 作为三角形方程组(5.2.12)的右端项,然后解出此三角形方程组,即可求得最小二乘解 x .

实际使用G-S过程时，由于 Q 的列是由 A 的列的线性组合产生的，一般说来，舍入误差较大，以致严重地影响 Q 阵的正交性，从而该算法的数值稳定性可能很差。为了提高算法的精度，将G-S过程进行适当修改，便得到所谓“修改的G-S正交化过程”。

算法5.2.1 修改的 Gram-Schmidt 方法(MGS 方法).
具体步骤如下.

1) 对 $k = 1, 2, \dots, n-1$, 做

$$1.1) \quad \mathbf{q}_k = \mathbf{a}_k / r_{kk}, \quad r_{kk} = \|\mathbf{a}_k\|, \quad (5.2.13)$$

1, 2) 对 $j = k + 1, \dots, n$ 做

$$1.2.1) \quad \mathbf{a}_j \leftarrow \mathbf{a}_j - r_{kj} \mathbf{q}_k, \quad r_{kj} = (\mathbf{q}_k, \mathbf{a}_j) \quad (5.2.14)$$

1.2.2) NEXT j ,

1.3) NEXT k .

$$2) \quad \text{置 } Q = (\mathbf{q}_1, \dots, \mathbf{q}_n), \quad R = (\mathbf{r}_1, \dots, \mathbf{r}_n), \quad (5.2.15)$$

$$\mathbf{r}_j^T = (r_{1j}, r_{2j}, \dots, r_{nj}, 0, \dots, 0).$$

$$3) \quad \text{解方程组 } R\mathbf{x} = Q^T \mathbf{b}. \quad (5.2.16)$$

$$4) \quad \text{输出 } \mathbf{x} = R^{-1} Q^T \mathbf{b}. \quad (5.2.17)$$

2.3 镜像变换法

我们知道，正交变换下任何向量的欧氏模是不变的。因此，如果 Q 是正交阵，那么

$$\|\mathbf{b} - A\mathbf{x}\| = \|Q(\mathbf{b} - A\mathbf{x})\| = \|Q\mathbf{b} - QA\mathbf{x}\|,$$

于是求 \mathbf{x} 使 $\|\mathbf{b} - A\mathbf{x}\|$ 达到最小的问题与求 \mathbf{x} 使 $\|Q\mathbf{b} - QA\mathbf{x}\|$ 达到最小的问题是等价的，其解均为最小二乘解。

根据第一章§3的理论，当 A 是列满秩矩阵时，存在正交阵 Q ，使

$$QA = \begin{pmatrix} R \\ 0 \end{pmatrix}, \quad (5.2.18)$$

其中 R 是 $n \times n$ 的上三角阵。令

$$Q\mathbf{b} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix} \begin{matrix} n \\ m-n \end{matrix},$$

于是

$$\begin{aligned} \|\mathbf{b} - A\mathbf{x}\| &= \left\| \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} - \begin{bmatrix} R \\ 0 \end{bmatrix} \mathbf{x} \right\| \\ &= [(\mathbf{b}_1 - R\mathbf{x})^T (\mathbf{b}_1 - R\mathbf{x}) + \mathbf{b}_2^T \mathbf{b}_2]^{\frac{1}{2}}. \end{aligned} \quad (5.2.19)$$

显然， $\mathbf{b}_1 - R\mathbf{x} = 0$ 的解便是使 (5.2.19) 达到最小的最小二乘

解。因此，不需要解法方程组(5.1.7)，而只需要解一次三角形方程组 $R\mathbf{x} = \mathbf{b}$ ，便可求出线性最小二乘解。

在 A 为列满秩阵的情况下，使用正交化方法或镜像变换法解“矛盾不太大的方程组”时，由于

$$\begin{aligned}\text{Cond}(R) &= \sqrt{\frac{R^T R \text{ 的最大特征值}}{R^T R \text{ 的最小特征值}}} \\ &= \sqrt{\frac{R^T Q^T Q R \text{ 的最大特征值}}{R^T Q^T Q R \text{ 的最小特征值}}} \\ &= \text{Cond}(QR) = \text{Cond}(A),\end{aligned}$$

即保持了条件数不变，从而避免了处理更为病态的方程组，因此用以上的方法去解最小二乘法问题，一般来说是可以获得比较满意的结果的。使用正交化方法和镜像变换法解矛盾方程组的具体算法，可以参看 Stewart 著的“矩阵计算引论”第五章。

§3 奇异值分解及广义逆矩阵

前面我们研究了在 $\text{null}(A) = 0$ 时求 $A\mathbf{x} = \mathbf{b}$ 的最小二乘解问题，当 $\text{null}(A) \neq 0$ 即矩阵 A 的秩 $r < n$ 时， $A\mathbf{x} = \mathbf{b}$ 的最小二乘解无论在解的表示形式和求解的方法上都将更为复杂，为了对后者进行较详细的研究，本节将引进**奇异值分解**和**广义逆矩阵**等概念。

3.1 奇异值分解

引理5.3.1 设矩阵 $A \in R^{m \times n}$ 的秩为 r ，则存在正交矩阵 $U \in R^{m \times m}$ ， $V \in R^{n \times n}$ ，使

$$U^T A V = \begin{bmatrix} \Delta & 0 \\ 0 & 0 \end{bmatrix}, \quad (5.3.1)$$

其中 Δ 是非奇异的 r 阶下三角阵。

证明 对于秩为 r 的矩阵 $A \in R^{m \times n}$ ，存在正交阵 Q 和排列

阵 P ，它们分别属于 $R^{m \times m}$ 和 $R^{n \times n}$ 使

$$QAP = \begin{bmatrix} \bar{A}_{11} & \bar{A}_{12} \\ 0 & 0 \end{bmatrix},$$

其中 \bar{A}_{11} 是 r 阶非奇异的上三角阵， \bar{A}_{12} 是 $r \times (n-r)$ 矩阵。有

$$(QAP)^T = \begin{bmatrix} \bar{A}_{11}^T & 0 \\ \bar{A}_{12}^T & 0 \end{bmatrix}.$$

对此又存在正交阵 $\bar{Q} \in R^{n \times n}$ ，使

$$\bar{Q}(QAP)^T = \begin{bmatrix} \Delta^T & 0 \\ 0 & 0 \end{bmatrix},$$

其中 Δ^T 是 r 阶非奇异的上三角阵，因此

$$(\bar{Q}^T(QAP)^T)^T = \begin{bmatrix} \Delta & 0 \\ 0 & 0 \end{bmatrix},$$

即

$$QAP\bar{Q}^T = \begin{bmatrix} \Delta & 0 \\ 0 & 0 \end{bmatrix}.$$

令 $Q^T = U$ ， $P\bar{Q}^T = V$ 。因为 Q ， \bar{Q} ， P 均为正交阵，所以 U ， V 也是正交阵，且分别属于 $R^{m \times m}$ 和 $R^{n \times n}$ ，使

$$U^TAV = \begin{bmatrix} \Delta & 0 \\ 0 & 0 \end{bmatrix},$$

其中 Δ 是 r 阶非奇异的下三角阵。

定理 5.3.1 设矩阵 $A \in R^{n \times n}$ 的秩为 r ，则存在 m 阶正交阵 U 和 n 阶正交阵 V ，使

$$U^TAV = D, \quad A = UDV^T, \quad (5.3.2)$$

$$D = \begin{bmatrix} D_r & 0 \\ 0 & 0 \end{bmatrix}, \quad D_r = \text{diag}(\mu_1, \mu_2, \dots, \mu_r),$$

其中 $\mu_i = \sqrt{\lambda_i}$ ，而 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$ 是 $A^T A$ 的非零特征值的全体。

证明 由引理5.3.1知道, 存在正交阵 \bar{U} 和 \bar{V} 使

$$\bar{U}^T A \bar{V} = \begin{bmatrix} \Delta & 0 \\ 0 & 0 \end{bmatrix},$$

其中 Δ 是 r 阶非奇异下三角阵. 因此 $\Delta^T \Delta$ 是 r 阶对称正定阵, 从而存在正交阵 V_r 使

$$\Delta^T \Delta = V_r \Lambda_r V_r^T$$

$$\Lambda_r = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_r \end{pmatrix},$$

这里 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$ 是 $\Delta^T \Delta$ 的全部特征值. 它们也是 $A^T A$ 的非零特征值. 令

$$D_r = \begin{pmatrix} \mu_1 & & \\ & \mu_2 & \\ & & \ddots \\ & & & \mu_r \end{pmatrix}, \quad \mu_i = \sqrt{\lambda_i} \quad (i = 1, 2, \dots, r)$$

则有 $D_r^2 = \Lambda_r$, 设

$$U_r = \Delta V_r D_r^{-1}, \quad (5.3.3)$$

则

$$\begin{aligned} U_r^T U_r &= (D_r^{-1})^T V_r^T \Delta^T \Delta V_r D_r^{-1} \\ &= D_r^{-1} V_r^T V_r \Lambda_r V_r^T V_r D_r^{-1} = D_r^{-1} \Lambda_r D_r^{-1} = I_r \end{aligned}$$

即 U_r 是正交阵, 从而

$$\bar{\bar{U}} = \begin{bmatrix} U_r & 0 \\ 0 & I_{m-r} \end{bmatrix}, \quad \bar{\bar{V}} = \begin{bmatrix} V_r & 0 \\ 0 & I_{n-r} \end{bmatrix}$$

分别为 m 和 n 阶正交阵. 又由(5.3.3)可得

$$U_r^T \Delta V_r = D_r$$

及
$$\bar{\bar{U}}^T \bar{U}^T A \bar{V} \bar{\bar{V}} = \bar{\bar{U}}^T \begin{bmatrix} \Delta & 0 \\ 0 & 0 \end{bmatrix} \bar{\bar{V}}$$

$$= \begin{bmatrix} U_r^T \Delta V_r & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} D_r & 0 \\ 0 & 0 \end{bmatrix}. \quad (5.3.4)$$

令 $\bar{\bar{U}} \bar{\bar{U}}^T = U$, $\bar{\bar{V}} \bar{\bar{V}}^T = V$, 则 U, V 都是正交阵, 且

$$U^T A V = \begin{bmatrix} D_r & 0 \\ 0 & 0 \end{bmatrix},$$

其中 $D_r = \text{diag}(\mu_1, \mu_2, \dots, \mu_r)$.

定义 5.3.1 设矩阵 $A \in R^{m \times n}$ 的秩为 r , $A^T A$ 的非零特征值为 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$, 则 $\mu_i = \sqrt{\lambda_i}$ ($i = 1, 2, \dots, r$) 称为 A 的奇异值, 等式 $A = U D V^T$ 称为 A 的奇异值分解.

在 A 的奇异值分解 $A = U D V^T$ 中, D 是由 A 的奇异值唯一确定的, 但是正交阵 U 和 V 的选取却不是唯一的. 当 $A^T A$ 有重特征值时, r 阶正交阵 V_r 的选取就不唯一, 从而 U 和 V 的选取也就不唯一. 当矩阵 A 是对称阵时, A 的奇异值就是它的特征值的绝对值.

有了引理 5.3.1 和定理 5.3.1 以后, 我们就可以在一般的情况下来讨论最小二乘法问题的解法了. 由于

$$\begin{aligned} \|b - Ax\| &= \|U^T b - U^T A x\| = \|U^T b - U^T A V V^T x\| \\ &= \|U^T b - \begin{bmatrix} \Delta & 0 \\ 0 & 0 \end{bmatrix} V^T x\|, \end{aligned} \quad (5.3.5)$$

$$\text{令 } c = U^T b = \begin{bmatrix} c_r \\ c_{n-r} \end{bmatrix}, \quad y = V^T x = \begin{bmatrix} y_r \\ y_{n-r} \end{bmatrix}, \quad (5.3.6)$$

将 (5.3.6) 代入 (5.3.5) 后得到

$$\begin{aligned} \|b - Ax\|^2 &= \left\| \begin{bmatrix} c_r \\ c_{n-r} \end{bmatrix} - \begin{bmatrix} \Delta & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} y_r \\ y_{n-r} \end{bmatrix} \right\|^2 \\ &= \left\| \begin{bmatrix} c_r - \Delta y_r \\ c_{n-r} \end{bmatrix} \right\|^2 = \|c_r - \Delta y_r\|^2 + \|c_{n-r}\|^2. \end{aligned}$$

因为 Δ 是非奇异的下三角阵, 所以方程组

$$\mathbf{c}_r - \Delta \mathbf{y}_r = 0 \quad (5.3.7)$$

有唯一解 $\bar{\mathbf{y}}_r$ ，令

$$\bar{\mathbf{x}} = \mathbf{V} \begin{bmatrix} \bar{\mathbf{y}}_r \\ \bar{\mathbf{y}}_{n-r} \end{bmatrix} = \begin{bmatrix} \bar{\mathbf{x}}_r \\ \bar{\mathbf{x}}_{n-r} \end{bmatrix}, \quad (5.3.8)$$

其中 $\bar{\mathbf{y}}_{n-r}$ 是任意 $(n-r)$ 维列向量。由于

$$\begin{aligned} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|^2 &\geq \|\mathbf{c}_{n-r}\|^2 = \|\mathbf{c}_r - \Delta \bar{\mathbf{y}}_r\|^2 + \|\mathbf{c}_{n-r}\|^2 \\ &= \|\mathbf{b} - \mathbf{A}\bar{\mathbf{x}}\|^2, \end{aligned}$$

故 $\bar{\mathbf{x}}$ 是 (5.1.2) 的最小二乘解。

如果取 $\bar{\mathbf{y}}_{n-r} = 0$ ，令

$$\bar{\mathbf{x}}^{(0)} = \mathbf{V} \begin{bmatrix} \bar{\mathbf{y}}_r \\ 0 \end{bmatrix}, \quad (5.3.9)$$

则

$$\|\bar{\mathbf{x}}\|^2 = \left\| \begin{bmatrix} \bar{\mathbf{x}}_r \\ \bar{\mathbf{x}}_{n-r} \end{bmatrix} \right\|^2 = \|\bar{\mathbf{x}}_r\|^2 + \|\bar{\mathbf{x}}_{n-r}\|^2 \geq \|\bar{\mathbf{x}}_r\|^2 = \|\bar{\mathbf{x}}^{(0)}\|^2,$$

所以 (5.3.9) 是最小二乘解中向量模为最小的解。

如果用 $\mathbf{A} = \mathbf{UDV}^T$ 代入 (5.3.5) 中，就得到

$$\|\mathbf{b} - \mathbf{A}\mathbf{x}\| = \|\mathbf{U}^T \mathbf{b} - \mathbf{U}^T \mathbf{UDV}^T \mathbf{x}\| = \|\mathbf{U}^T \mathbf{b} - \mathbf{DV}^T \mathbf{x}\|. \quad (5.3.10)$$

再将 (5.3.6) 代入 (5.3.10)，得到

$$\begin{aligned} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|^2 &= \left\| \begin{bmatrix} \mathbf{c}_r \\ \mathbf{c}_{n-r} \end{bmatrix} - \begin{bmatrix} \mathbf{D}_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y}_r \\ \mathbf{y}_{n-r} \end{bmatrix} \right\|^2 \\ &= \|\mathbf{c}_r - \mathbf{D}_r \mathbf{y}_r\|^2 + \|\mathbf{y}_{n-r}\|^2. \end{aligned}$$

类似于前面的讨论，可以得到 (5.1.2) 的向量模最小的最小二乘解

$$\begin{aligned} \bar{\mathbf{x}}^{(0)} &= \mathbf{V} \begin{bmatrix} \mathbf{D}_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{c}_r \\ \mathbf{c}_{n-r} \end{bmatrix} \\ &= \mathbf{V} \begin{bmatrix} \mathbf{D}_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} \mathbf{U}^T \mathbf{b}, \end{aligned} \quad (5.3.11)$$

从这个表达式中, 我们可以看到(5.1.2)的最小二乘解等于矩阵

$$V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T \quad (5.3.12)$$

与列向量 b 的乘积。而矩阵(5.3.12)与第二节中引入 \bar{A}^+ 有同样的作用, 我们也称之为 A 的广义逆矩阵。

3.2 广义逆矩阵

定义 5.3.2 对于矩阵 $A \in R^{m \times n}$, 当存在满足如下四个条件(称 Penrose 条件)的 $n \times m$ 矩阵 X 时, 称 X 为 A 的广义逆矩阵, 记为 A^+

$$AXA = A, \quad (p_1)$$

$$XAX = X, \quad (p_2)$$

$$AX = (AX)^T, \quad (p_3)$$

$$XA = (XA)^T, \quad (p_4)$$

根据广义逆矩阵的定义, 显然它具有如下性质:

$$1) \quad (A^+A)^2 = A^+A;$$

$$2) \quad (AA^+)^2 = AA^+.$$

定理 5.3.2 设矩阵 $A \in R^{m \times n}$ 的奇异值分解为

$$A = UDV^T, \quad D = \begin{bmatrix} D_r & 0 \\ 0 & 0 \end{bmatrix},$$

则

$$X = V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T$$

为 A 的广义逆矩阵 A^+ , 而且是唯一的。

证明 因为

$$\begin{aligned} AXA &= AV \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T A = UDV^T V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T UDV^T \\ &= UDV^T = A, \end{aligned}$$

$$\begin{aligned} XAX &= V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T U D V^T V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T \\ &= V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T = X, \end{aligned}$$

故

$$AX = U D V^T V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T = U \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix} U^T = (AX)^T.$$

类似可证 $XA = (XA)^T$.

所以 X 是 A 的广义逆矩阵 A^+ , 即

$$A^+ = V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T. \quad (5.3.13)$$

如果 X 和 \bar{X} 均为满足 $(p_1) \sim (p_4)$ 的广义逆矩阵, 则

$$\begin{aligned} X &= XAX = XA\bar{X}AX = X(AXA\bar{X})^T \\ &= X(A\bar{X})^T = XA\bar{X} = XA\bar{X}A\bar{X} = (XA)^T(\bar{X}A)^T\bar{X} \\ &= (\bar{X}AXA)^T = XAX = X. \end{aligned}$$

从而证明了 A^+ 的存在和唯一性.

从广义逆矩阵 A^+ 的唯一性知

$$A^+ = V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T,$$

它将不依赖于 A 的奇异值分解中的正交阵 U 和 V 的选取.

推理1 若矩阵 $A \in R^{m \times n}$ 的广义逆矩阵为 A^+ , 则 A^T 的广义逆矩阵为 $(A^+)^T$, 即

$$(A^+)^T = (A^T)^+.$$

证明 设 $A = UDV^T$, $D = \begin{bmatrix} D_r & 0 \\ 0 & 0 \end{bmatrix}$,

则

$$A^T = VDU^T = V \begin{bmatrix} D_r & 0 \\ 0 & 0 \end{bmatrix} U^T,$$

从而

$$\begin{aligned}(A^T)^+ &= (U^T)^T \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} V^T = \left(V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T \right)^T \\ &= (A^+)^T,\end{aligned}$$

即

$$(A^T)^+ = (A^+)^T.$$

推理2 $(A^+)^+ = A$.

其证明作为习题.

当 A 为列满秩阵时, 可以证明 §2 中定义的广义逆矩阵 $\bar{A} = (A^T A)^{-1} A^T$ 与本节定义的广义逆矩阵是一样的, 即 $\bar{A} = A^+$.

定理5.3.3 若 A 是列满秩阵, 则

$$(A^T A)^{-1} A^T = A^+;$$

若 A 是行满秩阵, 则

$$A^T (A A^T)^{-1} = A^+.$$

根据定义 5.3.2 和定理 5.3.2, 只要证明 $(A^T A)^{-1} A^T$ 满足 Penrose 条件就行了. 具体证明作为练习.

推理1 当 A 为列满秩阵时, 有 $A^+ A = I$; 当 A 为行满秩阵时, 有 $A A^+ = I$.

推理2 当 n 阶方阵 A 为满秩时, 则 $A^+ = A^{-1}$.

这个推论说明了 A^+ 是 A^{-1} 的推广.

定理 5.3.3 指出, 当 $A \in R^{m \times n}$ 的秩为 n 或 m 时, 广义逆矩阵可以表示成另一种形式. 现在我们将讨论当秩为 r 时, 如何用不同于 (5.3.13) 的形式来表示 A 的广义逆矩阵的问题. 先证明如下定理.

定理5.3.4 对于秩为 $r \neq 0$ 的矩阵 $A \in R^{m \times n}$, 必有列满秩阵 $B \in R^{m \times r}$ 与行满秩阵 $C \in R^{r \times n}$, 使

$$A = BC.$$

证明 对于秩为 $r \neq 0$ 的矩阵 A , 必有满秩矩阵 P 和 Q 使

$$PAQ = \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix}, \quad (5.3.14)$$

其中 I_r 是 r 阶单位矩阵. 设

$$P^{-1} = [B \ \bar{B}], \quad Q^{-1} = \begin{bmatrix} C \\ \bar{C} \end{bmatrix},$$

于是有

$$\begin{aligned} A &= P^{-1}(PAQ)Q^{-1} \\ &= [B \ \bar{B}] \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} C \\ \bar{C} \end{bmatrix} = BC. \end{aligned}$$

因为 P^{-1} 的 m 个列线性无关, 而 B 是由 P^{-1} 中的 r 个列组成, 所以 B 是列满秩阵, 同理 C 是行满秩阵.

定义 5.3.3 对秩为 $r \neq 0$ 的矩阵 $A \in R^{m \times n}$, 称分解式

$$A = BC \quad (5.3.15)$$

为 A 的**秩分解式**, 其中 B 和 C 分别是列满秩阵和行满秩阵.

定义 5.3.4 若 $A = BC$ 是非零矩阵 A 的秩分解式, 则

$$C^T(CC^T)^{-1}(B^TB)^{-1}B^T \quad (5.3.16)$$

叫做 A 的一个**广义逆矩阵**. 若 $A = O$, 则 O 就是 A 的广义逆矩阵. 按定理 5.3.3 的记法, (5.3.16) 可以写成

$$C^+B^+. \quad (5.3.17)$$

定理 5.3.5 A 的广义逆矩阵 C^+B^+ 满足 Penrose 条件 $(P_1) \sim (P_4)$.

证明 若记 $X = C^+B^+$ 则

$$AXA = BCC^T(CC^T)^{-1}(B^TB)^{-1}B^TBC = BC = A,$$

$$\begin{aligned} AX &= BCC^T(CC^T)^{-1}(B^TB)^{-1}B = B(B^TB)^{-1}B^T \\ &= (B(B^TB)^{-1}B^T)^T = (AX)^T. \end{aligned}$$

同理可证 $XAX = X$, $XA = (XA)^T$.

对于矩阵 A , 因为满足 Penrose 条件 $(P_1) \sim (P_4)$ 的广义逆矩阵是唯一的, 所以广义逆矩阵 C^+B^+ 可记为

$$A^+ = C^+B^+,$$

即

$$C^+B^+ = V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T.$$

例5.3.1 求

$$\begin{cases} x - y = -1, \\ x - y = 1, \\ 2x - y = 0 \end{cases}$$

的最小二乘解.

解 设

$$B = \begin{bmatrix} 1 & -1 \\ 1 & -1 \\ 2 & -1 \end{bmatrix},$$

显然 B 的秩为 2, 所以

$$B^+ = (B^T B)^{-1} B^T = \frac{1}{2} \begin{bmatrix} -1 & -1 & 2 \\ -2 & -2 & 2 \end{bmatrix},$$

于是

$$\begin{pmatrix} x \\ y \end{pmatrix} = B^+ \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} -1 & -1 & 2 \\ -2 & -2 & 2 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

例5.3.2 求 $Ax = b$ 的模为最小的最小二乘解, 其中

$$A = \begin{pmatrix} -1 & 2 & 1 \\ -1 & 2 & 1 \\ 0 & 3 & 2 \end{pmatrix}, \quad b = \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}.$$

解 显然 A 的秩为 2, 它的一个秩分解为

$$A = \begin{pmatrix} 1 & -1 \\ 1 & -1 \\ 2 & -1 \end{pmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 \\ 2 & -1 & 0 \end{bmatrix},$$

$$B^+ = (B^T B)^{-1} B^T = \frac{1}{2} \begin{bmatrix} -1 & -1 & 2 \\ -2 & -2 & 2 \end{bmatrix},$$

$$C^+ = C^T (CC^T)^{-1} = \frac{1}{14} \begin{pmatrix} 3 & 5 \\ 6 & -4 \\ 5 & -1 \end{pmatrix}$$

$$\begin{aligned} A^+ &= C^+ B^+ = \frac{1}{14} \begin{pmatrix} 3 & 5 \\ 6 & -4 \\ 5 & -1 \end{pmatrix} \frac{1}{2} \begin{bmatrix} -1 & -1 & 2 \\ -2 & -2 & 2 \end{bmatrix} \\ &= \frac{1}{28} \begin{pmatrix} -13 & -13 & 16 \\ 2 & 2 & 4 \\ -3 & -3 & 8 \end{pmatrix}, \end{aligned}$$

于是所求的解为

$$\mathbf{x} = A^+ \mathbf{b} = \frac{1}{28} \begin{pmatrix} -13 & -13 & 16 \\ 2 & 2 & 4 \\ -3 & -3 & 8 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix} = \frac{1}{28} \begin{pmatrix} 16 \\ 4 \\ 8 \end{pmatrix} = \frac{1}{7} \begin{pmatrix} 4 \\ 1 \\ 2 \end{pmatrix}.$$

下面我们用广义逆矩阵来给出方程组 $A\mathbf{x} = \mathbf{b}$ 的解及线性最小二乘解的一般表达式。

定理 5.3.6 线性方程组 $A\mathbf{x} = \mathbf{b}$ 有解的充要条件是

$$AA^+ \mathbf{b} = \mathbf{b} \quad (5.3.18)$$

当方程组有解时，它的通解是

$$\mathbf{x} = A^+ \mathbf{b} + (I - A^+ A) \boldsymbol{\xi}, \quad (5.3.19)$$

其中 $\boldsymbol{\xi}$ 是任意 n 维列向量。

当方程组无解时，(5.3.19) 是最小二乘法的通解（可理解为对应法方程组的通解），取 $\boldsymbol{\xi} = \mathbf{0}$ ，即得模为最小的最小二乘解 $A^+ \mathbf{b}$ 。

证明 1) 若方程组有解，即存在 \mathbf{y} 使 $A\mathbf{y} = \mathbf{b}$ ，由于 $A = AA^+A$ ，故

$$\mathbf{b} = \mathbf{A}\mathbf{y} = \mathbf{A}\mathbf{A}^+ \mathbf{A}\mathbf{y} = \mathbf{A}\mathbf{A}^+ \mathbf{b},$$

即 (5.3.18) 成立. 反之, 若 $\mathbf{A}\mathbf{A}^+ \mathbf{b} = \mathbf{b}$, 则 $\mathbf{A}^+ \mathbf{b} = \mathbf{x}$ 便是一个解, 从而证明了方程组有解的充要条件是 $\mathbf{A}\mathbf{A}^+ \mathbf{b} = \mathbf{b}$,

2) 当方程组有解时, 显然 $\mathbf{x}^{(0)} = \mathbf{A}^+ \mathbf{b}$ 是它的一个解. 现在要证齐次方程组 $\mathbf{A}\mathbf{x} = 0$ 的通解是

$$\bar{\mathbf{x}} = (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \boldsymbol{\xi}$$

其中 $\boldsymbol{\xi}$ 是任意 n 维列向量.

$$\text{因为 } \mathbf{A} \bar{\mathbf{x}} = \mathbf{A}(\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \boldsymbol{\xi} = (\mathbf{A} - \mathbf{A}\mathbf{A}^+ \mathbf{A}) \boldsymbol{\xi} = 0,$$

故 $\bar{\mathbf{x}} = (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \boldsymbol{\xi}$ 是 $\mathbf{A}\mathbf{x} = 0$ 的解, 另一方面, $\mathbf{A}\mathbf{x} = 0$ 的任何解 \mathbf{x} 可表为

$$\mathbf{x} = \mathbf{x} - \mathbf{A}^+ \mathbf{A}\mathbf{x} = (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \mathbf{x}.$$

又满足上式的 \mathbf{x} 必然满足 $\mathbf{A}\mathbf{x} = 0$, 从而证明了 $\mathbf{A}\mathbf{x} = 0$ 的通解是

$$\bar{\mathbf{x}} = (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \boldsymbol{\xi}$$

其中 $\boldsymbol{\xi}$ 是任意 n 维列向量, 综上所述, $\mathbf{A}\mathbf{x} = \mathbf{b}$ 的通解为

$$\mathbf{x} = \mathbf{A}^+ \mathbf{b} + (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \boldsymbol{\xi}.$$

3) 根据定理 5.1.1 和定理 5.1.2 可知, 最小二乘法必有解, 且其解必为法方程组

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$$

的解. 由 2) 知道, 该方程组的通解是

$$\mathbf{x} = (\mathbf{A}^T \mathbf{A})^+ \mathbf{A}^T \mathbf{b} + [\mathbf{I} - (\mathbf{A}^T \mathbf{A})^+ (\mathbf{A}^T \mathbf{A})] \boldsymbol{\xi},$$

其中 $\boldsymbol{\xi}$ 是任意 n 维列向量.

设 \mathbf{A} 的奇异值分解为 $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$, 则

$$\mathbf{A}^T \mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{U}^T \mathbf{U}\mathbf{D}\mathbf{V}^T = \mathbf{V}\mathbf{D}^2 \mathbf{V}^T$$

是 $A^T A$ 的奇异值分解, 所以

$$(A^T A)^+ = V \begin{bmatrix} (D_r^{-1})^2 & 0 \\ 0 & 0 \end{bmatrix} V^T,$$

于是

$$\begin{aligned} (A^T A)^+ A^T &= V \begin{bmatrix} (D_r^{-1})^2 & 0 \\ 0 & 0 \end{bmatrix} V^T V D U^T \\ &= V \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T = A^+, \end{aligned}$$

$$(A^T A)^- (A^T A) = ((A^T A)^+ A^T) A = A^+ A,$$

故法方程组的通解是

$$x = A^+ b + (I - A^+ A) \xi,$$

其中 ξ 是任意 n 维向量。当 $\xi = 0$ 时, $x = A^+ b$, 由广义逆矩阵的唯一性和 (5.3.11) 可知, $x = A^+ b$ 就是向量模为最小的最小二乘解。

§4 广义逆矩阵的一个递推算法

4.1 对称分块矩阵 A 的求逆

设对称分块矩阵为

$$A = \begin{bmatrix} P & R^T \\ R & Q \end{bmatrix},$$

其中 P 、 Q 也为对称矩阵且逆矩阵 P^{-1} 、 A^{-1} 存在。如果

$$\begin{bmatrix} P & R^T \\ R & Q \end{bmatrix} \begin{bmatrix} x_p \\ x_q \end{bmatrix} = \begin{bmatrix} b_p \\ b_q \end{bmatrix}, \quad (5.4.1)$$

则

$$\begin{bmatrix} x_p \\ x_q \end{bmatrix} = \begin{bmatrix} P & R^T \\ R & Q \end{bmatrix}^{-1} \begin{bmatrix} b_p \\ b_q \end{bmatrix}. \quad (5.4.2)$$

将(5.4.1)写成矩阵方程组的形式, 有

$$\begin{cases} P\mathbf{x}_p + R^T\mathbf{x}_q = \mathbf{b}_p, \end{cases} \quad (5.4.3)$$

$$\begin{cases} R\mathbf{x}_p + Q\mathbf{x}_q = \mathbf{b}_q. \end{cases} \quad (5.4.4)$$

类似于解二元一次方程组那样, 由(5.4.3)可得

$$\mathbf{x}_p = -P^{-1}(R^T\mathbf{x}_q - \mathbf{b}_p), \quad (5.4.5)$$

将(5.4.5)代入(5.4.4), 化简后可得

$$(Q - RP^{-1}R^T)\mathbf{x}_q = \mathbf{b}_q - RP^{-1}\mathbf{b}_p,$$

由于它的解存在且唯一, 所以

$$\mathbf{x}_q = -(Q - RP^{-1}R^T)^{-1}RP^{-1}\mathbf{b}_p + (Q - RP^{-1}R^T)^{-1}\mathbf{b}_q. \quad (5.4.6)$$

令 $Z = (Q - RP^{-1}R^T)^{-1}$, $Y = -ZRP^{-1}$, 代入(5.4.6), 得到

$$\mathbf{x}_q = Y\mathbf{b}_p + Z\mathbf{b}_q. \quad (5.4.7)$$

将(5.4.7)代入(5.4.5)化简后, 得到

$$\mathbf{x}_p = (P^{-1} - P^{-1}R^TY)\mathbf{b}_p + Y^T\mathbf{b}_q. \quad (5.4.8)$$

令 $W = P^{-1} - P^{-1}R^TY$, 代入(5.4.8)得到

$$\mathbf{x}_p = W\mathbf{b}_p + Y^T\mathbf{b}_q. \quad (5.4.9)$$

将(5.4.9)和(5.4.7)写成矩阵形式, 即得所求的解

$$\begin{bmatrix} \mathbf{x}_p \\ \mathbf{x}_q \end{bmatrix} = \begin{bmatrix} W & Y^T \\ Y & Z \end{bmatrix} \begin{bmatrix} \mathbf{b}_p \\ \mathbf{b}_q \end{bmatrix}. \quad (5.4.10)$$

比较(5.4.10)和(5.4.2)得到

$$\begin{bmatrix} P & R^T \\ R & Q \end{bmatrix}^{-1} = \begin{bmatrix} W & Y^T \\ Y & Z \end{bmatrix}, \quad (5.4.11)$$

其中

$$\begin{cases} Z = (Q - RP^{-1}R^T)^{-1}, \\ Y = -ZRP^{-1}, \\ W = P^{-1} - P^{-1}R^TY. \end{cases} \quad (5.4.12)$$

从 (5.4.12) 我们可以看出：分块求逆的方法是：先求出 P^{-1} ，再对 $(Q - RP^{-1}R^T)$ 求逆，然后利用矩阵的乘法求得 Y 和 W ，由此便可得到 A 的逆矩阵 A^{-1} 。

例 5.4.1 设

$$A = \begin{bmatrix} P & R^T \\ R & Q \end{bmatrix},$$

其中

$$P = Q = \begin{bmatrix} 3 & 0 & -1 & 0 \\ 0 & 3 & 0 & -1 \\ -1 & 0 & 3 & -1 \\ 0 & -1 & -1 & 3 \end{bmatrix}, \quad R = \begin{bmatrix} 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix},$$

求 A 的逆矩阵 A^{-1} 。

解 设

$$A^{-1} = \begin{bmatrix} W & Y^T \\ Y & Z \end{bmatrix},$$

由已知条件先求得

$$P^{-1} = \frac{1}{55} \begin{bmatrix} 21 & 1 & 8 & 3 \\ 1 & 21 & 3 & 8 \\ 8 & 3 & 24 & 9 \\ 3 & 8 & 9 & 24 \end{bmatrix}, \quad RP^{-1} = -\frac{1}{55} \begin{bmatrix} 8 & 3 & 24 & 9 \\ 3 & 8 & 9 & 24 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix},$$

$$RP^{-1}R^T = \frac{3}{55} \begin{bmatrix} 8 & 3 & 0 & 0 \\ 3 & 8 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

故

$$Z = (Q - RP^{-1}R^T)^{-1} = \frac{1}{2071} \begin{bmatrix} 963 & 127 & 377 & 168 \\ 127 & 963 & 168 & 377 \\ 377 & 168 & 939 & 369 \\ 168 & 377 & 369 & 939 \end{bmatrix};$$

$$Y = -ZRP^{-1} = \frac{1}{2071} \begin{pmatrix} 147 & 74 & 441 & 213 \\ 71 & 147 & 213 & 441 \\ 64 & 45 & 192 & 135 \\ 45 & 64 & 135 & 192 \end{pmatrix},$$

由Y便可求得 Y^T ;

$$W = P^{-1} - P^{-1}R^TY = \frac{1}{2071} \begin{pmatrix} 816 & 56 & 377 & 168 \\ 56 & 816 & 168 & 377 \\ 377 & 168 & 1131 & 504 \\ 168 & 377 & 504 & 1131 \end{pmatrix}.$$

由以上几项, A^{-1} 便可得知.

4.2 计算广义逆矩阵的递推算法

引理5.4.1 设

$$B = \begin{bmatrix} I_r \\ p_c^+ \end{bmatrix},$$

其中 p 是 n 维行向量, C^+ 是 $n \times r$ 阵,则

$$(BC)^+ = C^+B^+.$$

证明 设 $BC = A$, $C^+B^+ = X$, 由假设可得

$$B^TB = I_r + (PC^+)^T PC^+. \quad (5.4.13)$$

因为 B 是列满秩阵,所以 $B^+B = I$. 因此,

$$\begin{aligned} XA &= C^+B^+(BC) = C^+B^+(BC)(BC)^+BC = C^+C(BC)^+BC \\ &= C^T(C^+)^TC^TB^T((BC)^+)^T = C^TB^T((BC)^+)^T \\ &= [(BC)^T((BC)^+)^T]^T = (BC)^+(BC) = A^+A, \end{aligned} \quad (5.4.14)$$

于是

$$(XA)^T = (A^+A)^T = A^+A = XA, \quad AXA = AA^+A = A.$$

使用(5.4.13), 可以类似地证明

$$AX = A^+A, \quad (5.4.15)$$

于是

$$(AX)^T = (AA^+)^T = AA^+ = AX,$$

$$XAX = C^+B^+BCC^+B^+ = C^+CC^+B^+ = C^+B^+ = X.$$

因此, $X = C^+B^+$ 满足 Penrose 条件, 据定义 5.3.2, X 是 A 的广义逆矩阵, 即 $(BC)^+ = C^+B^+$.

引理 5.4.2 设

$$A = \begin{bmatrix} D & O \\ O & I \end{bmatrix} \begin{bmatrix} Q \\ p \end{bmatrix},$$

其中 D , Q 分别属于 $R^{t \times t}$ 和 $R^{t \times n}$ ($t < n$) 且 $DQ = C$ 为 C 的秩分解, p 为 n 维行向量并与 Q 的 t 行线性无关, 且 $d = p(I - C^+C) \neq 0$, 则

$$A^+ = (C^+ - d^+pC^+d^+). \quad (5.4.16)$$

证明 因为

$$A = \begin{bmatrix} D & O \\ O & I \end{bmatrix} \begin{bmatrix} Q \\ p \end{bmatrix}$$

是 A 的一个秩分解, 且

$$\begin{bmatrix} D & O \\ O & I \end{bmatrix}^+ = \begin{bmatrix} D^+ & O \\ O & I \end{bmatrix},$$

所以只需求

$$\begin{bmatrix} Q \\ p \end{bmatrix}^+.$$

据定理 5.3.3

$$\begin{aligned} \begin{bmatrix} Q \\ p \end{bmatrix}^+ &= [Q^T, p^T] \begin{bmatrix} QQ^T & Qp^T \\ pQ^T & pp^T \end{bmatrix}^{-1} = [Q^T, p^T] \begin{bmatrix} W & Y \\ Y^T & Z \end{bmatrix} \\ &= [Q^TW + p^Ty^T, Q^Ty + p^TZ]. \end{aligned} \quad (5.4.17)$$

由对称分块矩阵求逆法知道:

$$Z = [pp^T - pQ^T(QQ^T)^{-1}Qp^T]^{-1} = [pp^T - pQ^+Qp^T]^{-1},$$

$$Y^T = -ZpQ^T(QQ^T)^{-1} = -ZpQ^+,$$

$$W = (QQ^T)^{-1} - (QQ^T)^{-1}Qp^TY^T,$$

因为

$$C^+C = Q^+D^+DQ = Q^+Q,$$

而且

$$(C^+C)^T = C^+C,$$

$$\begin{aligned} dd^T &= p(I - C^+C)(I - C^+C)^T p^T = p(I - C^+C)p^T \\ &= pp^T - pC^+Cp^T = pp^T - pQ^+Qp^T, \end{aligned}$$

所以

$$Z = (dd^T)^{-1}. \quad (5.4.18)$$

因而有

$$\begin{aligned} Q^TY + p^TZ &= -Q^T(Q^+)^T p^T (dd^T)^{-1} + p^T (dd^T)^{-1} \\ &= (I - Q^+Q)^T p^T (dd^T)^{-1} = (I - C^+C)^T p^T (dd^T)^{-1} \\ &= d^T (dd^T)^{-1} = d^+, \end{aligned} \quad (5.4.19)$$

$$\begin{aligned} (Q^TW + p^TY^T)D^+ &= [Q^T(QQ^T)^{-1} - Q^T(QQ^T)^{-1}Qp^TY^T \\ &\quad - p^TZpQ^+]D^+ = [Q^+ + Q^+Qp^T(dd^T)^{-1}pQ^+ - p^T(dd^T)^{-1}pQ^+]D^+ \\ &= Q^+D^+ + Q^+Qp^T(dd^T)^{-1}pQ^+D^+ - p^T(dd^T)^{-1}pQ^+D^+ \\ &= C^+ + C^+Cp^T(dd^T)^{-1}pC^+ - p^T(dd^T)^{-1}pC^+ \\ &= C^+ - (I - C^+C)p^T(dd^T)^{-1}pC^+ \\ &= C^+ - d^T(dd^T)^{-1}pC^+ = C^+ - d^+pC^+. \end{aligned} \quad (5.4.20)$$

由(5.4.17)、(5.4.19)、(5.4.20)可得

$$\begin{aligned} A^+ &= \begin{bmatrix} Q \\ p \end{bmatrix}^+ \begin{bmatrix} D^+ & 0 \\ 0 & 1 \end{bmatrix} = [Q^TW + p^TY^TQ^TY + p^TZ] \begin{bmatrix} D^+ & 0 \\ 0 & 1 \end{bmatrix} \\ &= [(Q^TW + p^TY^T)D^+Q^TY + p^TZ] \\ &= (C^+ - d^+pC^+d^+). \end{aligned}$$

引理5.4.3 设 I 是 n 阶单位阵, $x = (x_1, x_2, \dots, x_n)$ 为 n 维行向量, 则

$$(I + x^Tx)^{-1}x^T = x^T / (1 + \|x\|^2). \quad (5.4.21)$$

证明 因为据Householder变换 Q , 可将任一行向量 x^T 变成

$$Qx^T = [x \ e,$$

其中 e 为单位向量。因此,

$$\begin{aligned}
 Q(I + x^T x)^{-1} x^T &= Q(I + x^T x)^{-1} Q^T Q x^T \\
 &= [Q(I + x^T x)^{-1} Q^T] Q x^T = [Q(I + x^T x) Q^T]^{-1} Q x^T \\
 &= [I + (Q x^T)(Q x^T)^T]^{-1} Q x^T \\
 &= [I + \|x\|^2 e e^T]^{-1} \|x\| e \\
 &= \frac{\|x\|}{1 + \|x\|^2} e = \frac{Q x^T}{1 + \|x\|^2}, \quad (5.4.22)
 \end{aligned}$$

在(5.4.22)的两边乘 Q^T , 即得

$$(I + x^T x)^{-1} x^T = x^T / (1 + \|x\|^2).$$

定理5.4.1 设 $C \in R^{k \times n}$, $p \in R^{1 \times n}$, 且

$$A = \begin{bmatrix} C \\ p \end{bmatrix},$$

则

$$A^+ = [C^+ - f p C^+ f], \quad (5.4.23)$$

$$f = \begin{cases} d^+ & \text{当 } d \neq 0 \text{ 时;} \\ C^+ (C^T)^+ p^T h & \text{当 } d = 0 \text{ 时;} \end{cases}$$

$$d = p(I - C^+ C), \quad h = (1 + p C^+ (C^T)^+ p^T)^{-1}.$$

证明 当 $d \neq 0$ 时, 设 $C = DQ$ 为 C 的秩分解。因为 $t+1 \leq n$, 否则 $C^+ C = I$ 与 $d \neq 0$ 矛盾; 而且 p 与 Q 的行向量线性无关, 否则, 必有 $a \in R^{1 \times t}$ 使 $p = aQ$, 于是

$$p C^+ C = p (Q^+ D^+ D Q) = p (Q^+ Q) = a Q Q^+ Q = a Q = p, \text{ 这也和}$$

$d \neq 0$ 矛盾, 因此 $\begin{bmatrix} Q \\ p \end{bmatrix}$ 是行满秩阵, 由此得到了 A 的一个秩分解

$$A = \begin{bmatrix} C \\ p \end{bmatrix} = \begin{bmatrix} D & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} Q \\ p \end{bmatrix}.$$

根据引理5.4.2得到

$$A^+ = \begin{bmatrix} Q \\ p \end{bmatrix}^+ \begin{bmatrix} D^+ & 0 \\ 0 & 1 \end{bmatrix} = (C^+ - d^+ p C^+ d^+). \quad (5.4.24)$$

当 $d = o$ 时, $p - pC^+C = o$, 即 $p = pC^+C$, 故

$$A = \begin{bmatrix} C \\ p \end{bmatrix} = \begin{bmatrix} C \\ pC^+C \end{bmatrix} = \begin{bmatrix} I \\ pC^+ \end{bmatrix} C = BC.$$

设

$$B^+ = \begin{bmatrix} I \\ pC^+ \end{bmatrix}^+ = (W, u).$$

并作下三角阵

$$(B, v) = \begin{bmatrix} I & O \\ pC^+ & 1 \end{bmatrix}, \quad (5.4.25)$$

其中 $v^T = (O^T, 1)$, 显然 v^T 与 B^T 的行线性无关, 且

$$\begin{bmatrix} B^T \\ v^T \end{bmatrix} = \begin{bmatrix} I & O \\ O & 1 \end{bmatrix} \begin{bmatrix} B^T \\ v^T \end{bmatrix},$$

由引理5.4.2证明的结果可以得到

$$\begin{bmatrix} B^T \\ v^T \end{bmatrix}^+ = [(B^T)^+ - (d_1^T)^+ v^T (B^T)^+ (d_1^T)^+]$$

或

$$[Bv]^+ = \begin{bmatrix} B^+ - B^+ v d_1^+ \\ d_1^+ \end{bmatrix}, \quad (5.4.26)$$

其中

$$\begin{aligned} d_1 &= (I - BB^+)v = \left[I - \begin{bmatrix} I \\ pC^+ \end{bmatrix} (W, u) \right] \begin{bmatrix} O \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} -u \\ 1 - pC^+u \end{bmatrix}. \end{aligned} \quad (5.4.27)$$

又由于

$$(B, v)^+ = (B, v)^{-1} = \begin{bmatrix} I & O \\ -pC^+ & 1 \end{bmatrix}, \quad (5.4.28)$$

比较(5.4.26)与(5.4.28)可得

$$d_1^+ = (-pC^+, 1). \quad (5.4.29)$$

又由(5.4.27)可得

$$\begin{aligned} d_1^T d_1 &= v^T (I - BB^+) v = v^T d_1 \\ &= (0, 1) \begin{bmatrix} -u \\ 1 - pC^+ u \end{bmatrix} = 1 - pC^+ u. \end{aligned}$$

故

$$d_1^+ = (d_1^T d_1)^{-1} d_1^T = (-(1 - pC^+ u)^{-1} u^T, 1). \quad (5.4.30)$$

比较(5.4.30)与(5.4.29)得到

$$pC^+ = (1 - pC^+ u)^{-1} u^T. \quad (5.4.31)$$

解(5.4.31)得

$$\begin{aligned} (pC^+)^T - (pC^+)^T (pC^+ u) &= u, \\ (pC^+)^T &= (I + (pC^+)^T pC^+) u, \\ u &= [I + (pC^+)^T pC^+]^{-1} (pC^+)^T, \end{aligned} \quad (5.4.32)$$

因为 pC^+ 是 n 维行向量, 根据引理 5.4.3, (5.4.32) 可以写成

$$u = [1 + pC^+ (pC^+)^T]^{-1} (pC^+)^T. \quad (5.4.33)$$

又由于 B 为列满秩阵, 据定理 5.3.3 的推论 1 可得

$$I = B^+ B = (W, u) \begin{bmatrix} I \\ pC^+ \end{bmatrix} = W + u pC^+,$$

即

$$W = I - u pC^+. \quad (5.4.34)$$

由于

$$A = BC$$

且

$$B = \begin{bmatrix} I \\ pC^+ \end{bmatrix}.$$

据引理 5.4.1 可得

$$\begin{aligned} A^+ &= C^+ B^+ = C^+ (W, u) = C^+ (I - u pC^+, u) \\ &= (C^+ - C^+ u pC^+, C^+ u). \end{aligned}$$

令

$$\begin{aligned} f &= C^+ u = C^+ [1 + pC^+ (pC^+)^T]^{-1} (pC^+)^T \\ &= C^+ (pC^+)^T [1 + pC^+ (pC^+)^T]^{-1} = C^+ (pC^+)^T h, \end{aligned} \quad (5.4.35)$$

即

$$A^+ = (C^+ - f p C^+, f).$$

根据定理5.4.1, 矩阵

$$A = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_m \end{pmatrix} \quad (\text{其中 } \alpha_i \text{ 为 } n \text{ 维行向量})$$

的广义逆矩阵 A^+ , 可以使用如下递推算法.

算法5.4.1 本算法可以用依次计算 $A_1 = \alpha_1$.

$A_2 = \begin{bmatrix} A_1 \\ \alpha_2 \end{bmatrix}, \dots, A_m = \begin{bmatrix} A_{m-1} \\ \alpha_m \end{bmatrix}$ 的广义逆矩阵的方法, 最后算得 A

的广义逆矩阵 A^+ , 其具体步骤如下.

1) 置 A_0^+ , A_0 分别是 $R^{n \times 1}$ 和 $R^{1 \times n}$ 中的零阵;

2) 对 $k = 1, 2, \dots, m$,

2.1) 置 $b_k = \alpha_k A_{k-1}^+$

2.1.1) 置 $d = \alpha_k - b_k A_{k-1}$,

2.2) 如果 $d \neq 0$, 转2.4),

2.3) 置 $h = (1 + b_k b_k^T)^{-1}$,

2.3.1) 置 $f = A_{k-1}^+ b_k^T h$,

2.4) 置 $f = d^T (d d^T)^{-1}$,

2.5) 置 $A_k^+ = [A_{k-1}^+ - f b_k, f]$,

2.6) NEXT k .

3) 输出结果 A_m^+ .

本算法略去了计算过程中的存储技巧和计算细节.

例5.4.1 应用递推算法计算

$$C = \begin{pmatrix} -1 & 2 & 1 \\ -1 & 2 & 1 \\ 0 & 3 & 2 \end{pmatrix}$$

的广义逆矩阵 C^+ 。

解 设 $C_1 = [-1, 2, 1]$

则

$$c_1^+ = c_1^T (c_1 c_1^T)^{-1} = [-1, 2, 1]^T \cdot \frac{1}{6} = \frac{1}{6} \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}.$$

设

$$C_2 = \begin{bmatrix} c_1 \\ -1 & 2 & 1 \end{bmatrix}, \quad \text{求 } C_2^+.$$

由于

$$\begin{aligned} d &= [-1, 2, 1] \left[I - \frac{1}{6} \begin{pmatrix} 1 & -2 & -1 \\ -2 & 4 & 2 \\ -1 & 2 & 1 \end{pmatrix} \right] \\ &= [-1, 2, 1] \begin{pmatrix} \frac{5}{6} & \frac{2}{6} & \frac{1}{6} \\ \frac{2}{6} & \frac{2}{6} & -\frac{2}{6} \\ \frac{1}{6} & -\frac{2}{6} & \frac{5}{6} \end{pmatrix} = [0, 0, 0], \end{aligned}$$

所以, 取

$$\begin{aligned} f &= c_1^+ (c_1^+)^T p^T h = c_1^+ (c_1^+)^T p^T [1 + p c_1^+ (c_1^+)^T p^T]^{-1} \\ &= c_1^+ (p c_1^+)^T [1 + p c_1^+ (p c_1^+)^T]^{-1} = \frac{1}{6} \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix} \cdot \frac{1}{2} = \frac{1}{12} \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}, \\ C_2^+ &= \left(\frac{1}{6} \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix} - \frac{1}{12} \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}, \frac{1}{12} \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix} \right) = \frac{1}{12} \begin{pmatrix} -1 & -1 \\ 2 & 2 \\ 1 & 1 \end{pmatrix}, \end{aligned}$$

设 $C_3 = \begin{bmatrix} C_2 \\ 0 \ 3 \ 2 \end{bmatrix}$, 求 C_3^+ .

由于

$$\begin{aligned} d &= (0, 3, 2) \left(I - \frac{1}{12} \begin{pmatrix} -1 & -1 \\ 2 & 2 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} -1 & 2 & 1 \\ -1 & 2 & 1 \end{pmatrix} \right) \\ &= (0, 3, 2) \begin{pmatrix} \frac{5}{6} & \frac{1}{3} & \frac{1}{6} \\ \frac{1}{3} & \frac{1}{3} & -\frac{1}{3} \\ \frac{1}{6} & -\frac{1}{3} & \frac{5}{6} \end{pmatrix} = \left(\frac{4}{3}, \frac{1}{3}, \frac{2}{3} \right) \neq 0, \end{aligned}$$

所以, 取

$$f = d^+ = \begin{pmatrix} \frac{4}{3} \\ \frac{1}{3} \\ \frac{2}{3} \end{pmatrix} \left(\begin{pmatrix} \frac{4}{3}, \frac{1}{3}, \frac{2}{3} \end{pmatrix} \begin{pmatrix} \frac{4}{3} \\ \frac{1}{3} \\ \frac{2}{3} \end{pmatrix} \right)^{-1} = \frac{1}{7} \begin{pmatrix} 4 \\ 1 \\ 2 \end{pmatrix},$$

$$\begin{aligned} C^+ &= C_3^+ = (C_3^+ - f p C_3^+, f) = \begin{pmatrix} \frac{1}{12} \begin{pmatrix} -1 & -1 \\ 2 & 2 \\ 1 & 1 \end{pmatrix} - \frac{2}{21} \begin{pmatrix} 4 & 4 \\ 1 & 1 \\ 2 & 2 \end{pmatrix} \frac{1}{7} \begin{pmatrix} 4 \\ 1 \\ 2 \end{pmatrix} \\ \frac{1}{28} \begin{pmatrix} -13 & -13 & 16 \\ 2 & 2 & 4 \\ -3 & -3 & 8 \end{pmatrix} \end{pmatrix}. \end{aligned}$$

线性最小二乘法问题, 曾经一度只用法方程组 $A^T A x = A^T b$ 来求解. 但是, 由于方程组的条件数之间存在着 $\text{Cond}(A^T A) = (\text{Cond}(A))^2$ 关系, 所以当 $\text{Cond}(A)$ 较大时, 法方程组就变得非常“病态”, 用这种方法便难于求得满意的解了. 此时, 为了

避免处理非常病态的方程组，可以采用正交化方法和镜像变换法，“直接”对矛盾方程组进行求解。

还须注意，当矩阵 A 的秩小于 n 时，最小二乘解不唯一，从而需要求模为最小的最小二乘解。为此，需要用到奇异值分解，这种分解的难点是，确定 (5.3.2) 中哪些 u_i 应视为 0。Golub 与 Reinsch 于 1970 年提出了一种数值上十分稳定的分解方法，再加上奇异值的稳定性（见习题 2），意味着这种方法是比较可靠的。我们可以把这样得到的解用广义逆矩阵来表示。目前广义逆矩阵已有很多计算方法，它可以不通过奇异值分解而直接进行计算，由于篇幅有限，在此，我们仅介绍了一种递推算法。

第五章 习 题

5.1 证明 (1) $(AA^+)^2 = AA^+$; (2) $(A^+A)^2 = A^+A$; (3) $(A^+)^+ = A$.

5.2 设 $\sigma_i(A)$ 和 $\sigma_i(A+E)$ 分别为矩阵 A 和 $A+E$ 的奇异值，证明 $\sigma_i(A) - \|E\| \leq \sigma_i(A+E) \leq \sigma_i(A) + \|E\|$.

5.3 给出应用奇异值分解求解齐次线性方程组的方法。

5.4 证明

若 A 是列满秩阵，则 $(A^T A)^{-1} A^T = A^+$;

若 A 是行满秩阵，则 $A^T (A A^T)^{-1} = A^+$.

5.5 证明定理 5.3.5 中的等式 $X = X A X$ 和 $X A = (X A)^T$ 成立。

5.6 证明等式 (5.4.15)。

5.7 设 $A \in R^{m \times n}$ 是列满秩阵， $C \in R^{l \times n}$ 是行满秩阵，导出一个在约束条件 $Cx = d$ 下，使 $\|b - Ax\|$ 达到最小的最小二乘解的算法，其中 b 和 d 均为已知的列向量。

5.8 求

$$\begin{cases} y_1 + y_2 + y_3 = 1, \\ 2y_1 - y_2 = 0 \end{cases}$$

的模为最小的解。

5.9 证明 在 A 的奇异值分解 $A = UDV^T$ 中, V 的行和 U 的列分别是 $A^T A$ 和 AA^T 的特征向量。

5.10 利用对称分块矩阵求逆法, 求

$$A = \begin{pmatrix} 3 & 0 & -1 & 0 \\ 0 & 3 & 0 & -1 \\ -1 & 0 & 3 & -1 \\ 0 & -1 & -1 & 3 \end{pmatrix}$$

的逆矩阵 A^{-1} 。

5.11 矛盾方程组 $Ax = b$ 的最小二乘解(5.3.8)与通解(5.3.19)是否等价? 为什么?

5.12 利用逆推算法计算

$$A = \begin{pmatrix} 1 & 3 & 1 & 1 \\ 2 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

的广义逆矩阵 A^+ 。

5.13 给定数据表

x_i	0.000	1.445	2.890	4.335	5.780
y_i	1.8419	2.9633	18.2360	98.7410	529.2178

求形如 $y = ae^{bx}$ 函数的最小二乘方拟合。

5.14 求 $Ax = b$ 的模为最小的最小二乘解, 其中

$$A = \begin{pmatrix} 2 & 1 & -1 \\ 0 & 3 & -2 \\ 2 & 4 & -3 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix}.$$

参 考 书

- [1] Stewart. G. "Introduction to Matrix Computation" Academic Press. New York, 1973.
- [2] 曹志浩等编, 矩阵计算和求方程根。
- [3] 党诵诗编, 矩阵论及其在测绘中的应用, 测绘出版社, 1980.

第六章 稀疏矩阵

有限元方法、结构分析、网络分析及数学规划中的许多问题，都将归结为求解大型线性稀疏方程组的问题：

$$Ax = b. \quad (6.0.1)$$

这种方程组的特点是：矩阵 A 的阶数 n 很大，但其非零元素的个数 τ 所占的比例却很小。若以 τ/n^2 表示矩阵 A 的稀疏度，通常当 τ/n^2 小于 $0.05 \sim 0.25$ 时，即认为 A 是稀疏的。

稀疏矩阵技术，就是针对矩阵的稀疏性来建立相应算法的一种技术。其基本出发点是，只让非零的信息参与运算和存贮，这样往往能在存贮量和机器时间的节省带来巨大的收益。限于篇幅，本章只着重讨论用直接法求解稀疏线性方程组的问题。

§1 稀疏矩阵的存贮

1.1 等带宽矩阵的存贮

对于矩阵 A ，如果存在正整数 m ，使

$$a_{ij} = 0, \quad |i - j| > m \quad (i, j = 1, \dots, n,) \quad (6.1.1)$$

则称 A 是**等带宽**的，称 $2m + 1$ 为其带宽， m 为半带宽，把满足 $|i - j| \leq m$ 的所有的 (i, j) 元素的集合叫做 A 的**带形区域**。

当对方程 (6.0.1) 进行不带行交换的高斯消去法时，在带形区域之外将不会产生非零元素。这时，我们只需把 A 的带形区域部分的元素存入机器。

注意到矩阵 A 的带形区域的形状，见图 6.1 的阴影部分，它只需补充少数几个元素就可当作二维数组来存放。

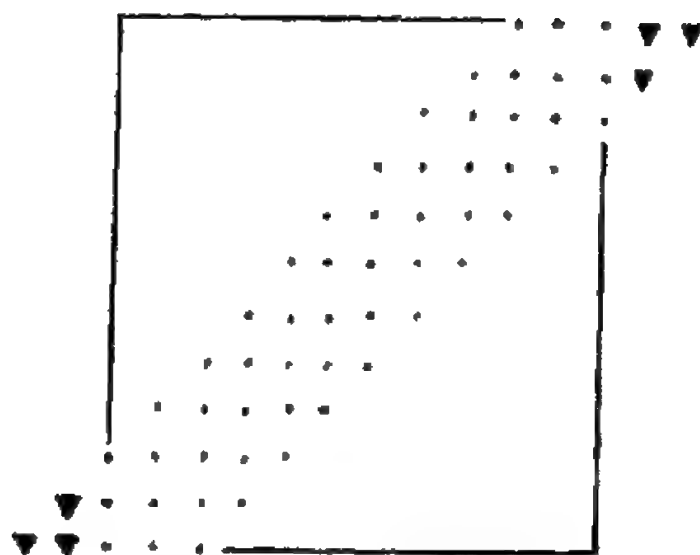


图6.1 等带宽矩阵

图中画有“▼”的元素，其数值可以任意给定。把得到的二维数组记为 $L[1:n, 1:2m+1]$ ；用 i, j 表示元素 a_{ij} 在 A 中的行、列，用 i', j' 表示 a_{ij} 在 L 中的行列，则

$$\begin{aligned} i' &= i, \quad j' = j - i + m + 1, \\ i &= 1, 2, \dots, n, \\ \max(1, i - m) &\leq j \leq \min(n, i + m). \end{aligned} \quad (6.1.2)$$

(6.1.2)指出了 A 在带形区域中任意一个元素的位置与该元素在 L 中的位置的检索关系。所谓矩阵的压缩存贮，主要就是要建立这样的位置检索关系。当 A 是对称阵时，等带宽矩阵可以只存贮其下三角阵部分，设将 a_{ij} 存在二维数组 $L[1:n, 1:m+1]$ 的 (i, j) 处，则

$$\begin{aligned} i' &= i, \quad j' = i - j + 1, \quad j \leq i, \\ i &= 1, 2, \dots, n, \quad j = 1, \dots, \max(1, i - m). \end{aligned}$$

1.2 变带宽矩阵的存贮

如果 A 是对称阵，我们把 A 的第 i 行上对角线左侧离对角线最远的非零元素到同行对角元间的距离记为 β_i ，称为第 i 行

上的半带宽。当 A 是正定对称阵时，由于可以进行 **Cholesky** 分解

$$A = LL^T,$$

其中 L 是下三角阵，可以证明， L 阵在第 i 行上的半带宽一定与 β_i 相等。这样，当把 A 的三角分解的输出 L 的元素放在 A 的下三角阵中对应的元素上时，则 A 和 L 的存贮量要求将是完全一致的。

A 的变带宽存贮是这样实现的：逐行把 A 的下三角部分的每一行的元素（包括半带宽范围内的每一个零元和非零元），依次存入一个一维数组 $L[1:s]$ 中，其中

$$s = \sum_{i=1}^n (\beta_i + 1), \quad (6.1.3)$$

$$0 \leq \beta_i \leq i-1, \quad i=1, \dots, n.$$

令 A 的第 i 个对角元在 L 数组中的地址为 $d[i]$ ，则 $a_{ij} (j \leq i)$ 在 L 中的地址 l 可表为

$$l = d[i] + j - 1, \quad (6.1.4)$$

$$d[i] = \sum_{k=1}^i (\beta_k + 1). \quad (6.1.5)$$

反之，由 L 数组中的第 l 个数，也可以求出该数在原矩阵 A 中的行号 i 和列号 j ：

$$i = \max_k (l > d[k]) + 1,$$

$$j = l - d[i-1] + i - \beta_i = l - d[i] + 1. \quad (6.1.6)$$

可见，在组织对称阵 A 的变带宽输入时，应输入两个一维数组 $L[1:s]$ 和 $D[1:n]$ ，前者是逐行存放 A 在各条半带宽上的元素的，后者是存放 A 的对角元在 L 中的位置的。

如图 6.2，设各有关数据为： $n=9$ ， $\beta_1=0$ ， $\beta_2=1$ ， $\beta_3=2$ ， $\beta_4=1$ ， $\beta_5=0$ ， $\beta_6=3$ ， $\beta_7=2$ ， $\beta_8=2$ ， $\beta_9=3$ ； $s=23$ ， $d_1=1$ ， $d_2=3$ ，

$d_3 = 6, d_4 = 8, d_5 = 9, d_6 = 13, d_7 = 16, d_8 = 19, d_9 = 23$; a_{76} 在 L 中的地址为: $l = d_7 + 6 - 7 = 15$, L 中的第 12 个数的行号 i 、列号 j 各为:

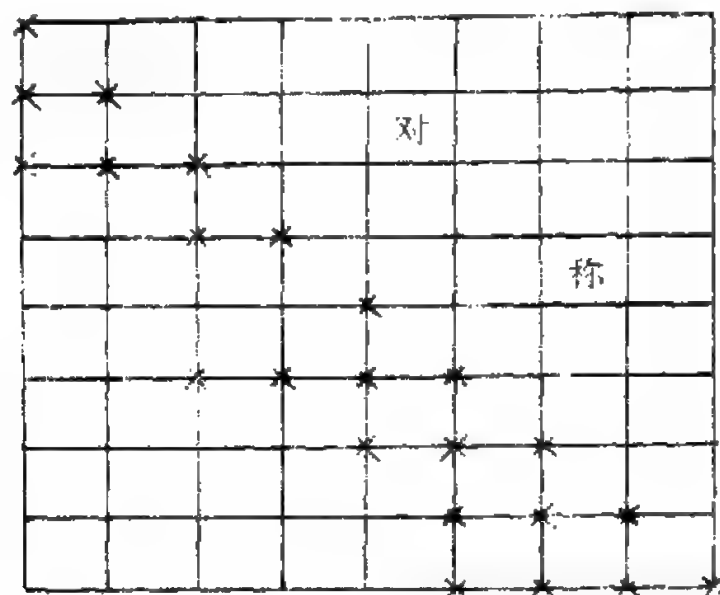


图6.2 对称变带宽矩阵

$$i = \max_k (12 > d[k]) + 1 = 5 + 1 = 6.$$

$$j = 12 - d_6 + 6 = 5.$$

1.3 压缩存贮法

当矩阵 A 中的非零元素的分布十分零乱时, 可用下述的压缩存贮法进行存贮。

在内存中指定 n 个连续编号的单元, 称为列首址表, 记为 $BC[1:n]$, $BC[j]$ 的内容被规定为 A 的第 j 列第 1 个非零元素的信息条首址表所谓一个非零元素 a_{ij} 的信息条, 是指由三个连续编号的单元组成的成组单元, 三个连续单元中的内容按图 6.3 规定。

信息条首址	第二地址	第三地址
元素行号 i	元素 a_{ij}	下一非零元信息条首址或零

图6.3 信息条内容

如图，如果在第 j 列的 a_{ij} 元素之下还有后续的非零元，就在第三地址的内容中填写下一个非零元的信息条首址，如果 a_{ij} 是第 j 列中最后一个非零元，就在第三地址中写零。这样，当知道了一个元素的信息条首址之后，不仅同时能确定该元素的值，而且连下一个非零元的信息条首址也知道了。由于同列的非零元的所有的信息条是一一连接的，所以这种存贮方法又有连接表法之称。

当知道元素的行号 i 和列号 j 时，查找元素 a_{ij} 的步骤应当是：

1) 由 j 找 $BC[j]$ ，设其内容为：

$$(BC[j]) = d,$$

d 是第 j 列第 1 个非零元的信息首址。

2) 如果 $(d) = i$ 则 $(d+1) = a_{ij}$ ，于是检索 a_{ij} 的工作结束；

3) 如果 $(d) \neq i$ ，则 $(d+2) \Rightarrow d$ ，并转 2)；只要元素的行、列信息没有填错，总有一次会使 2) 成立。

设矩阵 A 的非零元的总数为 τ ，则存贮全部信息所需的存贮单元总数为 $3\tau + n$ ，当 $3\tau + n \ll n^2$ 时，这对存贮的节约是很有利的。

当用直接法求解方程组 (6.0.1) 时，原矩阵中原来是零的元素，在消元过程中可能变成非零的元素，原来是非零的元素又可能消元为零，于是在信息条的“链条”中，必须讨论“增添”和“删除”一个元素的问题。上述存贮方法能方便地实现这种要求。

首先，假定在内存中有一个有足够容量的空白单元区，其首址为 L ，它是由三个连续编号的单元组成的空白条的“链条”组成的，空白条与信息条的区别在于：空白条的第一、第二地址的内容可以是任意的，但第三地址的内容必须填写下一个空

白条的首址。而空白区的首址 L 存放在一个固定的地址 S 中，即令

$$(S) = L.$$

于是当要在矩阵的第 j 列中增加一个元素时，例如，设在第 j 列中，相邻的两个非零元 $a_{i_1,j}$ 、 $a_{i_2,j}$ 已经分别放在首址为 d_1 、 d_2 的相连接的信息条之中，其中 $i_2 - i_1 \geq 2$ ，

d_1	$d_1 + 1$	$d_1 + 2$...	d_2	$d_2 + 1$	$d_2 + 2$...
i_1	$a_{i_1,j}$	d_2		i_2	$a_{i_2,j}$	p	

今欲在第 j 列中的 $a_{i_1,j}$ 、 $a_{i_2,j}$ 之间插入一元素 $a_{i,j}$ ($i_1 < i < i_2$)，应按下列步骤进行：

1) 求 $(s) = L$ ，找到空白区首址；

2) $(L+2) \Rightarrow s$ ，把下一空白条首址作为空白区首址；

3) 在 $L, L+1, L+2$ 中，分别记入： $i, a_{i,j}, d_2$ ，即“征用”原空白区的第 1 个空白条为信息条以记录 $a_{i,j}$ 及其行号，并把新信息条的下端接好；

4) 把 L 记入 $d_1 + 2$ ，即把新信息条的上端也接好。

当 $a_{i,j}$ 是第 j 列中的第 1 个要增新的非零元或最末一个要增新的非零元时，亦可仿此进行连接。

当要从第 j 列中“删除”一个元素时，例如，设 $a_{i_1,j}$ 、 $a_{i_2,j}$ 、 $a_{i_3,j}$ 已经分别存放在以 d_1 、 d_2 、 d_3 为首址的互相连接的信息条之中，即

d_1	$d_1 + 1$	$d_2 + 2$...	d_2	$d_2 + 1$	$d_2 + 2$...	d_3	$d_3 + 1$	$d_3 + 2$
i_1	$a_{i_1,j}$	d_2		i_2	$a_{i_2,j}$	d_3		i_3	$a_{i_3,j}$	p

今欲删除以 d_2 为首址的信息条，为使被删除的信息条能“退役”到空白区中，以备将来征用，可按下列步骤实现：

1) $(d_2 + 2) \Rightarrow d_1 + 2$, 这就完成了对 d_2 条的删除,

2) 由 $(s) = L$, 求出 L , 并令 $d_2 \Rightarrow s$, $L \Rightarrow d_2 + 2$, 即把删除掉的信息条作为空白区的第 1 个空白条, 并使之与原有空白条连接起来。

还可设计出一些压缩存贮的方法, 但使用上述连接表方法, 却独具一种优点, 即增删一个元素都很方便, 特别, 欲变换矩阵的两列时, 例如, 欲交换第 j_1 和 j_2 两列, 只需把 $B[j_1]$ 和 $B[j_2]$ 的内容互相交换即可。这种思想还可很容易地推广到复稀疏结构的情形。

§2 消元次序

2.1 引言

在稀疏矩阵计算中, 在设计或选用一种算法时, 应尽可能使计算过程所需的最大存贮量和运算次数达于极小, 为了达到这个目标, 先介绍几个术语。

填充量 如图 6.4

所示的 7 阶矩阵中, 用 \times 表示非零元素, 空白处为零元素。在进行消元过程的第一步时, 应把第 1 列上三角化, 此时第 1 列上的 $(2, 1)$,

$$\begin{bmatrix} \times & & \times & & \times & & \\ \times & \times & \otimes & & \otimes & \times & \\ \times & & \times & & \otimes & & \\ & & & \times & & & \\ & & \times & & \times & & \\ & & & & & \times & \\ \times & & \otimes & & \otimes & \times & \end{bmatrix}$$

图 6.4 填充量

$(3, 1)$, $(7, 1)$ 应消元成零。与此同时, 原来有些是零的元素如 $(2, 3)$, $(2, 6)$, $(3, 6)$, $(7, 3)$, $(7, 6)$ 却将由零变成非零。图中用 \otimes 标出的地方, 就是这些新增的非零的元素。在本例中, 其个数是 5 个, 称为第一步消元过程中产生的填充量。如以此

数减去原第 1 列中的被消元成零的非零元个数，本例中应为 $5 - 3 = 2$ ，可称之为第 1 次消元时的填充增量。类似地，在第 2 次消元时也会出现新的填充量增量，而在整个消元过程中，填充量增量将有一个最大的积累。

值得指出的是：填充量的积累与主元选取的次序或矩阵的行列的排列次序有关。如图 6.5，当按自然次序消元时，则在第一步消元时，填充量即达于最大；但若把第 1 行、列与最末的行列交换时，则整个消元过程中，填充量保持为零。下一节我们将介绍一些重排行、列次序的方法，以使填充量积累趋于极小。

$$\begin{array}{cc}
 \left[\begin{array}{ccccc}
 \times & \times & \times & \times & \times \\
 \times & \times & & & \\
 \times & & \times & & \\
 \times & & & \times & \\
 \times & & & & \times
 \end{array} \right] & \left[\begin{array}{ccccc}
 \times & & & & \times \\
 & \times & & & \times \\
 & & \times & & \times \\
 & & & \times & \times \\
 \times & \times & \times & \times & \times
 \end{array} \right] \\
 (a) & (b)
 \end{array}$$

图 6.5 填充量与行列交换

主元容限 在利用列主元或全主元消去法把矩阵上三角化的过程中，例如，在做第 1 步列主元消去时，通常应把第 1 列中绝对值最大的元素换行到 (1,1) 的位置，然后进行消元。这样做的目的，一般是为了提高计算精度的。但是，如果当前的主元的选取恰好导致填充量增加最大，这就不利于内存的节省，于是宁愿选择一个绝对值不是最大但又不会引起填充量有过大增长的元素作主元。当然，作为主元的元素，其绝对值也不能过小，否则就会引起精度的大的丧失。人为地规定一个界限 $\epsilon > 0$ ，当矩阵的元素依模大于此数时，该元素就具有了当主元的资格，如果它引起的填充量也果然不大，就可把主元定下

来。所规定的这个界限 ϵ 即称为主元容限。这个限量可以经验地给以规定，但它应体现计算精度及存贮节省的统一要求。

主元容限的可供推荐的选取方法是：先求

$$\alpha = \max_{i,j} |a_{ij}|$$

并取 $\epsilon = 10^{-4}\alpha$ 。

主元容限的上述选取方法是以矩阵具有一定的均衡性为前提的，即矩阵的元素依模均不大于 1，且每个行和每个列中都至少有一个元素依模不小于 1/2。当矩阵非奇异时，通常总能找到两个对角阵 D_1 和 D_2 ，使 $D_1 A D_2$ 能成为均衡或接近均衡。

2.2 填充量极小化方法

在进行某一步消元时，该步中的填充量极小化问题是容易处理的。然而，即使每一步消元都保证了该步中的填充量极小化，也未必能保证得到全局的填充量极小化方案。要得到全局的使填充量积累最小的消元次序，原则上应比较所有可能的消元次序方案中实际填充量积累的大小，才能选出最好的方案。但是，所有可能的消元次序方案总数可达 $(n!)^2$ 之多！因此，现实的作法仍只需保证每步消元过程中的填充量极小就行了。但此种作法只能叫作局部的填充量极小化方法。

由于每一步消元过程中填充量极小化的讨论方法是一样的，所以我们只讨论第 1 步消元时的填充量极小化方法。

令 B 是 A 的布尔矩阵，其元素规定为

$$b_{ij} = \begin{cases} 1, & \text{当 } a_{ij} \neq 0 \text{ 时;} \\ 0, & \text{当 } a_{ij} = 0 \text{ 时。} \end{cases} \quad (6.2.1)$$

设在第 1 步消元过程中，以 a_{pq} 为主元。则 $a_{pq} \neq 0$ ，称 p 为主行， q 为主列。于是，在 (i, l) 元素处产生填充量的充要条件是

$$(a_{pi} \neq 0) \wedge (a_{lq} \neq 0) \wedge (a_{il} = 0) \quad (6.2.2)$$

或

$$(b_{ij} = 1) \wedge (b_{iq} = 1) \wedge (b_{ij} = 0) \quad (6.2.3)$$

或

$$\mathbf{e}_p^T \mathbf{B} \mathbf{e}_i \mathbf{e}_i^T \bar{\mathbf{B}} \mathbf{e}_i \mathbf{e}_i^T \mathbf{B} \mathbf{e}_q = 1. \quad (6.2.4)$$

这里 $\bar{\mathbf{B}}$ 是 \mathbf{B} 的共轭矩阵, 它是把 \mathbf{B} 中的 0 换成 1, 1 换 0 后所成的矩阵. 因此, 当以 a_{pq} 为主元时, 除 p 行和 q 列上的元素外, 在其余各处产生填充量的总和为

$$g_{r,q} = \sum_{\substack{i,j \\ i \neq p, j \neq q}} \mathbf{e}_p^T \mathbf{B} \mathbf{e}_i \mathbf{e}_i^T \bar{\mathbf{B}} \mathbf{e}_i \mathbf{e}_i^T \mathbf{B} \mathbf{e}_q. \quad (6.2.5)$$

注意, 在上式中, 当 $i = p$ 或 $j = q$ 时, 由于

$$\mathbf{e}_p^T \mathbf{B} \mathbf{e}_i \mathbf{e}_i^T \bar{\mathbf{B}} \mathbf{e}_i = 0 \text{ 及 } \mathbf{e}_i^T \bar{\mathbf{B}} \mathbf{e}_q \mathbf{e}_i^T \mathbf{B} \mathbf{e}_q = 0,$$

故 (6.2.5) 中求和号下的 $i \neq p, j \neq q$ 的限制可以去掉, 得

$$\begin{aligned} g_{r,q} &= \sum_{i,j=1}^n \mathbf{e}_p^T \mathbf{B} \mathbf{e}_i \mathbf{e}_i^T \bar{\mathbf{B}} \mathbf{e}_i \mathbf{e}_i^T \mathbf{B} \mathbf{e}_q \\ &= \sum_{i,j=1}^n \mathbf{e}_p^T \mathbf{B} \mathbf{e}_i \mathbf{e}_i^T \bar{\mathbf{B}}^T \mathbf{e}_i \mathbf{e}_i^T \mathbf{B} \mathbf{e}_q \\ &= \mathbf{e}_p^T \mathbf{B} \left(\sum_{i=1}^n \mathbf{e}_i \mathbf{e}_i^T \right) \bar{\mathbf{B}}^T \left(\sum_{i=1}^n \mathbf{e}_i \mathbf{e}_i^T \right) \mathbf{B} \mathbf{e}_q \\ &= \mathbf{e}_p^T \mathbf{B} \bar{\mathbf{B}}^T \mathbf{B} \mathbf{e}_q. \end{aligned} \quad (6.2.6)$$

令

$$g_{\alpha\beta} = \min_{\substack{p,q \\ |a_{pq}| > \epsilon}} \{ \mathbf{e}_p^T \mathbf{B} \bar{\mathbf{B}}^T \mathbf{B} \mathbf{e}_q \}, \quad (6.2.7)$$

则 (α, β) 应是第 1 次消元时产生的使填充量极小化的主元位置. 这个方法是 1967 年由 Tewarson 提出的, 其后, Markowitz 又把此法加以简化. 称

$$b_{pj}b_{iq}$$

为第 1 步消元时, 如以 (p, q) 为主元, 而在 (i, j) 处可能产生的填充量, 由于

$$b_{pj}b_{iq} \geq b_{ij}(1 - b_{ii})b_{iq}$$

类似于(6.2.6)的推导, 可得

$$\begin{aligned}\bar{g}_{p,q} &= \sum_{\substack{i \neq p \\ j \neq q}} b_{ij} b_{ij} = \mathbf{e}_p^T (\mathbf{B} - \mathbf{I}) \mathbf{V} \mathbf{V}^T (\mathbf{B} - \mathbf{I}) \mathbf{e}_q \\ &= \mathbf{e}_p^T (\mathbf{B} - \mathbf{I}) \mathbf{M} (\mathbf{B} - \mathbf{I}) \mathbf{e}_q,\end{aligned}\quad (6.2.8)$$

这里用 \mathbf{V} 表示分量全为1的向量, \mathbf{M} 表示元素全为1的矩阵。

令

$$\bar{g}_{\alpha,\beta} = \min_{\substack{p,q \\ |a_{p,q}| \geq \epsilon}} \{ \mathbf{e}_p^T (\mathbf{B} - \mathbf{I}) \mathbf{M} (\mathbf{B} - \mathbf{I}) \mathbf{e}_q \}, \quad (6.2.9)$$

这就得到了 Markowitz 的填充量极小化的策略。

注意, 由于(6.2.8)中 $\mathbf{e}_p^T (\mathbf{B} - \mathbf{I}) \mathbf{V}$ 和 $\mathbf{V}^T (\mathbf{B} - \mathbf{I}) \mathbf{e}_q$ 分别表示与元素 $a_{p,q}$ 同行或同列的不包含 $a_{p,q}$ 的非零元的总数, 因此 $\bar{g}_{p,q}$ 的计算十分容易。在实践中, 此法已被广泛的采用。

当矩阵 A 为对称阵时, 为了保证消元过程中流动矩阵 $A^{(k)}$ 的对称性, 通常应把主元选在对角线上, 即于(6.2.9)中, 应取 $\alpha = \beta$, 得

$$\bar{g}_\alpha = \min_{\substack{p \\ |a_{pp}| \geq \epsilon}} \{ \mathbf{V}^T (\mathbf{B} - \mathbf{I}) \mathbf{e}_p \}^2. \quad (6.2.10)$$

2.3 计算量的极小化方法

在消元过程中, 欲使总的计算量趋于极小, 选取主元时也应有一定的策略: 我们仍以第1次消元时的计算量分析为主。

当 $a_{p,q}$ 被选为主元时, 其进行过程为: 先作一次除法, 以求主元的倒数 $1/a_{p,q}$, $\mathbf{V}^T \mathbf{B} \mathbf{e}_q - 1$ 次乘法以实现主元除主列, 这是矩阵三角分解中需要存贮的信息; 再作 $\mathbf{e}_p^T \mathbf{B} \mathbf{V} - 1$ 次乘法, 实现主元除主行, 这是消元时所必须的信息; 最后作

$$\sum_{\substack{i \neq p \\ j \neq q}} (\mathbf{e}_p^T \mathbf{B} \mathbf{e}_j) (\mathbf{e}_i^T \mathbf{B} \mathbf{e}_q) = (\mathbf{e}_p^T \mathbf{B} \mathbf{V} - 1) (\mathbf{V}^T \mathbf{B} \mathbf{e}_q - 1)$$

次乘法, 这是消元时需要作的。于是, 1次消元时的乘除法总量为:

$$\begin{aligned}\bar{g}_{pq} &= 1 + (\mathbf{V}^T \mathbf{B} \mathbf{e}_p - 1) + (\mathbf{e}_p^T \mathbf{B} \mathbf{V} - 1) + (\mathbf{e}_p^T \mathbf{B} \mathbf{V} - 1) \\ &\quad \times (\mathbf{V}^T \mathbf{B} \mathbf{e}_q - 1) = \mathbf{e}_p^T \mathbf{B} \mathbf{M} \mathbf{B} \mathbf{e}_q,\end{aligned}\quad (6.2.11)$$

比较(6.2.8)与(6.2.11), 并引入

$$\begin{aligned}h_{pq}(\mu) &= \mathbf{e}_p^T (\mathbf{B} - \mu \mathbf{I}) \mathbf{M} (\mathbf{B} - \mu \mathbf{I}) \mathbf{e}_q, \\ (0 \leq \mu \leq 1)\end{aligned}\quad (6.2.12)$$

则有

$$h_{pq}(0) = \bar{g}_{pq}, \quad h_{pq}(1) = \bar{g}_{pq}.\quad (6.2.13)$$

上式分别表示, 当 a_{pq} 为主元时, 欲使计算量极小应取 $\mu = 0$, 而欲使可能的填充量极小应取 $\mu = 1$. 于是 $h_{pq}(\mu)$ 应视为一次消元过程中, 当 a_{pq} 为主元时, 使计算量和填充量达于极小的带权平均值. 而

$$\bar{h}_{\alpha, \beta}(\mu) = \min_{\substack{p, q \\ |a_{pq}| \geq \varepsilon}} h_{pq}(\mu) \quad (6.2.14)$$

将是综合计算量的极小化和填充量的极小化的选取主元的策略. 权因子的选取, 取决于机器的硬件和软件, 当计算速度成为计算过程中的主要矛盾时, 宜取 μ 接近于零, 反之, 则应取 μ 接近于1.

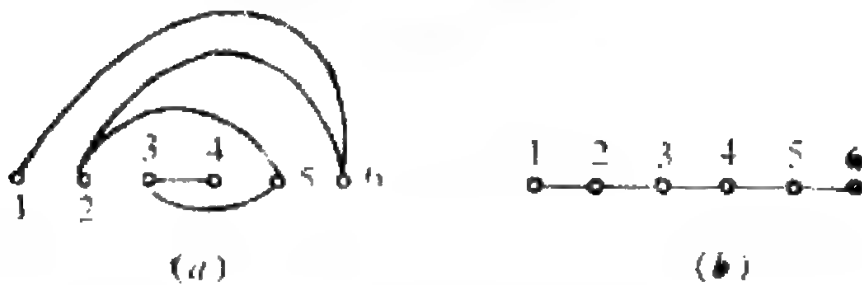
2.4 带宽极小化方法

由于对称带状阵在带形区域内的零元在存贮时也要分配单元, 当带宽较大而带形区域内的零元又较多时, 常可利用交换行的相似变换使带宽减小. 如图6.6(a), 对给定的对称6阶矩阵, 由于半带宽 $m = 6$, 依等带宽的压缩形式存贮时需占用21个单元. 但若将原矩阵的第1, 2, 3, 4, 5, 6各行和各列依次排在第1, 6, 2, 5, 3, 4各行和各列时则得如图6.6(b)中的稀疏结构, 此时半带宽为 $m = 1$, 仅需占用11个单元.

这就自然地提出下述问题: 求排列阵 P , 使

$$PAP^T$$

的带宽最小，其中 A 是对称正定矩阵。



我们将借用图论中的一些术语来讨论上述问题。

先设 A 为任意的 n 阶实矩阵，并于平面上任取互异的 n 个顶点 v_1, v_2, \dots, v_n 。当矩阵 A 的元素 $a_{ij} \neq 0$ 时，用一条边把顶点 v_i 和 v_j 连接起来，记为 $\langle v_i, v_j \rangle$ ，称为边。当视 $\langle v_i, v_j \rangle$ 与 $\langle v_j, v_i \rangle$ 为同一条边时，称这种边为 **无向边**，否则称为 **有向边**， $\langle v_i, v_i \rangle$ 称为 **自身路**，若记顶点的集合为 $V = \{v_1, v_2, \dots, v_n\}$ ，边的集合记为 $E = \{\langle v_i, v_j \rangle\}$ ，则顶点与边的集合

$$G = \{V, E\} \quad (6.2.15)$$

称为图。

当所有的边均是元向边时，相应的图叫做**无向图**，反之叫做**有向图**。当 A 是对称阵时，通常我们只研究它的无向图。例如，前面已给的图6.6 (a), (b) 中的矩阵的无向图为 (见图6.7)。

如果在图 G 和图 G' 的顶点和边之间可以建立一一对应的关系，并且对应边上的对应点也是对应的，则称图 G 与图 G'

同形，同形的图可视为同一个图。当 A 是对称阵时，显然其同形的无向图是唯一确定的，反之，同形的无向图也唯一地确定了对称阵的非零结构。当 A 是正定对称阵时，由于矩阵 A 的每个对角元均是正数，因此，其无向图上的每个点都有一条自身路，为了简化图面，今后在讨论正定对称阵的无向图时，我们将略去其所有的自身路。

交换对称阵的第 i 、 j 两行和两列，比较矩阵在交换前后的无向图，可知，只要在原有的无向图中，交换顶点 v_i 和 v_j 的编号，即可得到交换后的矩阵的无向图；由此推知，矩阵

$$PAP^T \quad (6.2.16)$$

乃是，在 A 的无向图上重新对顶点进行编号后的无向图的矩阵。于是求排列阵的问题就转变成如何对图的顶点进行编号，使 PAP^T 的带宽极小化的问题。

在图 G 中任意取定一个顶点 v_i ，称通过此点的边数为该点的度，这些边中不包括 v_i 本身的顶点叫做 v_i 的邻点，显然 v_i 的邻点的个数与 v_i 的度是相等的。

注意到矩阵 A 的第 i 行上的半带宽 β_i ，乃是 v_i 的编号减去 v_i 的邻点中比 v_i 编号小的最小的编号，当 v_i 的邻点中没有比 v_i 更小的编号时，则 $\beta_i = 0$ 。可见，寻求顶点的编号的方法，以使图的相应的矩阵有最小的带宽，其原则之一应当是，每个顶点 v 的编号与 v 的邻点的编号不宜相差过大。

Cuthill 和 Makee 基于图的理论，用对顶点重新编号的方法，使矩阵的带宽减小，文献上称为 CM 方法，其后又为 Rosen 加以改进，发展成为 RCM 方法。目前，它已应用于许多大型结构的程序分析系统。

于图 $G = \{V, E\}$ 的顶点的集合 V 中，任取度数最小的一个顶点 v 称为根，记

$$L_1 = \{v\} \quad (6.2.17)$$

称为**第一分层**，把 V 中与 v 相邻的点的集合（不包括 v ）记为 L_2 ，称为第 2 个分层，一般地把与 L_{i-1} 相邻但不属于 L_{i-1} 与 L_{i-2} 的点的集合记为 L_i ，叫做第 i 个分层；设 V 中的点一共可区分为 k 个分层，称

$$L_1, L_2, \dots, L_k \quad (6.2.18)$$

为图的以 v 为根的一个分层结构。显然

$$L_i \cap L_j = \emptyset \ (i \neq j), \quad \bigcup_{i=1}^k L_i = V \quad (6.2.19)$$

称每个分层 L_i 中的点数 w_i 为 L_i 的宽度，而把

$$w = \max_{1 \leq i \leq k} w_i \quad (6.2.20)$$

叫做**分层结构的宽度**， k 为**分层结构的深度**。根的选取，可能不是唯一的，选定以后，即将它编号成 1，然后依次给第二层中的点编号，在给第 i 层中的顶点编号时，应先编第 $i-1$ 层中度数最低的那些点的邻点。如这样的点不止一个，则应使第 i 层中编号大的点与编号小的点分别与第 $i-1$ 层编号大的点和编号小的点相邻。当各层均已编完后，一个排列阵 P 即已形成。以图 6.7(a) 为例，如以 v_1 为根，则按 CM 方法，各分层分别为：

$$L_1 = \{v_1\}, L_2 = \{v_6\}, L_3 = \{v_2\}, L_4 = \{v_5\}, L_5 = \{v_3\}, L_6 = \{v_4\},$$

因为原来的点的编号 1, 2, 3, 4, 5, 6 应重排成 1, 3, 5, 6, 4, 2，即新的点的编号应为 1, 6, 2, 5, 3, 4，于是

$$P = (e_1, e_6, e_2, e_5, e_3, e_4).$$

当以 P, P^T 左乘、右乘于图 6.6(a) 中的矩阵时，即得图 6.6(b) 中所示的矩阵。

给图的顶点重新排序，使相应矩阵的带宽趋于极小，这是排序过程中所要追求的一个目标，此外，使矩阵的外形趋于极

小，则是排序中所要追求的另一个目标。

矩阵的外形是指

$$\beta = \sum_{i=1}^n (\beta_i + 1), \quad (6.2.21)$$

其中 β_i 是矩阵的第 i 行上的半带宽。显然，当用变带宽方式存贮矩阵时，矩阵的外形的减小将有利于我们的讨论。

在 CM 方法的基础上如果使用反编号的方法，即把由 CM 方法得到的某顶点的编号 i 重新编号为 $n-i+1$ ，可达到使外形减小的目的。例如，以图 6.8 为图的矩阵，其外形为 109，而用反排序方法即 RCM 方法排序时，外形化为 97。

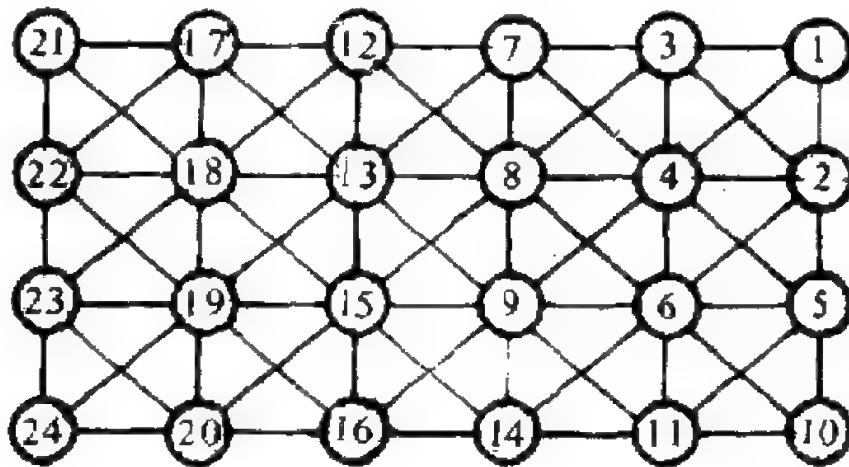


图6.8 用 CM 方法排序

RCM 方法也还存在着可以改进的地方。首先，RCM 方法需对所有度数较低的点进行搜索，逐个计算分层结构，每次计算中都包括求一个排列矩阵 P 和 $PA P^T$ 的带宽，这些工作量是比较大的。为了减少这种搜索的次数，Gibbs, Poole, Stochmeyes 等人基于一种几何直觉，提出一种伪直径概念，设计了另一种求分层结构的方法，称为 GPS 方法。在某些情况下，使用这种方法确实显露了它的优点。

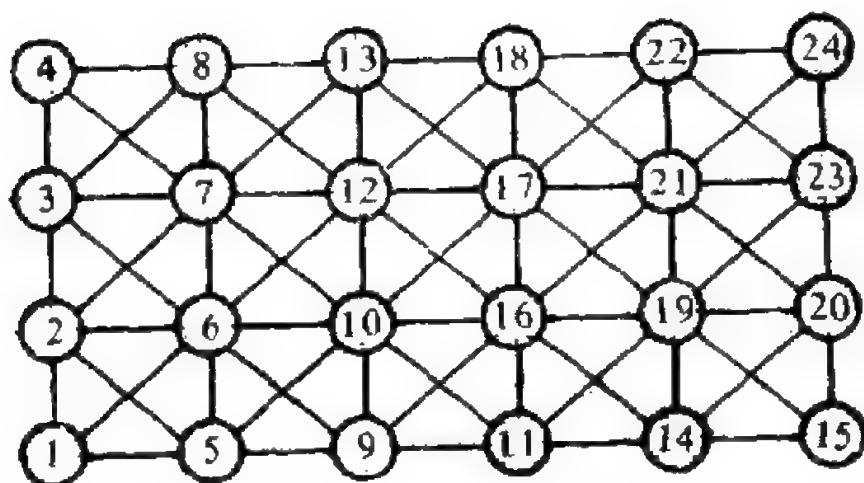


图 6.9 用 RCM 方法排序

几何直觉的启示是，分层结构的深度越大，则各分层结构的平均宽度将越小，注意 L_i 分层上一个点 v 的编号 r 减去 L_{i-1} 分层上与 v 点相邻且编号小于 r 的最小编号的号码就是矩阵 A 的第 i 行上的半带宽，这个值的上界必定不超过 L_i, L_{i-1} 的宽度之和，所以分层结构的深度大，矩阵的带宽将平均地越小。设以 v 为根的分层结构是图的具有最大深度 k 的一个分层结构，我们把从 v 到第 k 个分层上任意一点 u 之间的最短路径叫做图的一个直径，其长为 k ，而 u, v 两点即称为直径的端点。显然，要求一个图的一条直径及其端点，理论上应进行大量的搜索，这是不现实的，简单而实用的作法是：

- 1) 取任意一个度数最低的点 v ；
- 2) 以 v 为根做分层结构： L_1, L_2, \dots, L_k ，
- 3) 在 L_k 上依度数增加的次序选点 s ，以 s 为根做分层结构，如果深度不增加且 s 已经选完则转4)，如未选完，继续做3)，如深度有增加，则 s 送到 v 并转2)；

4) L_k 上任意一点与 v 之间的最短路径都是 k ，可当作是图的直径的近似值，特称之为伪直径。于 L_k 上找一点 u ，使以 u 为根做的分层结构所对应的矩阵的带宽 $w(u)$ 满足

$$w(u) = \min_{s \in L} w(s). \quad (6.2.22)$$

例如，在图(6.8)中，取 $v = 1$ ，则分层结构的各分层为：

$$L_1 = \{1\}, L_2 = \{2, 3, 4\}, L_3 = \{5, 6, 7, 8, 9\},$$

$$L_4 = \{10, 11, 12, 13, 14, 15, 16\}, L_5 = \{17, 18, 19, 20\},$$

$$L_6 = \{21, 22, 23, 24\};$$

这时 $k = 6$ ，由于 $w(21) = w(24) = 7$ ， $w(22) = w(23) = 6$ ，取 $u = 23$ ，以 $u = 23$ 为根的分层结构的层为

$$M_1 = \{23\}, M_2 = \{24, 20, 19, 18, 22\},$$

$$M_3 = \{16, 15, 13, 12, 17, 21\}, M_4 = \{14, 9, 8, 7\},$$

$$M_5 = \{11, 6, 4, 3\}, M_6 = \{10, 5, 2, 1\}$$

反编 M_i 的层号，即记

$$M'_i = M_{k+1-i}, \quad i = 1, 2, \dots, k,$$

由于图中每个顶点都将分别属于分层结构 L 和分层结构 M' 的一个分层，因此，每个点各对应 L 分层的一个足码 i 和 M' 的一个足码 j ，亦即对应一个足码对 (i, j) ，如图(6.10)，这个足码对的表称为图 6.8 的以 1, 23 为伪直径端点的相关层表示。如果把相关层表示图中，具有相关足码为 (i, i) 的那些点组成另一分层结构 N 的第 i 层，则此时

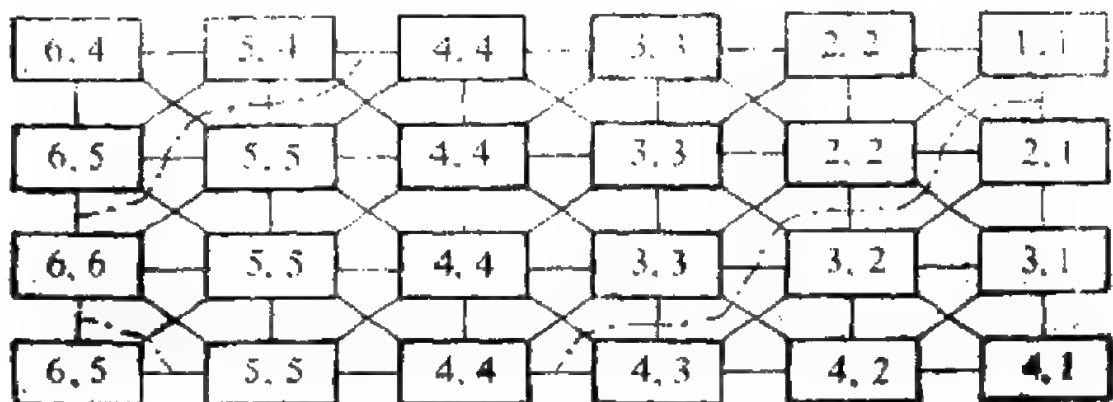


图6.10 图6.8的相关层表示

$$N_1 = \{1\}, N_2 = \{3, 4\}, N_3 = \{7, 8, 9\},$$

$$N_4 = \{12, 13, 15, 16\}, N_5 = \{20, 19, 18\}, N_6 = \{23\}.$$

它们表明，当以某一伪直径的端点为根时， N_i 中的点都应在同一层次。剩下还需把相关足码为 $(i, j) i \neq j$ 的点分别编入 N 的各层中，这些点在相关层表示图中，一般将分隔成互相隔离的几个连通的子集：

$$c_1, c_2, \dots, c_t.$$

在上例中， $t=3$ 而：

$$c_1 = \{2, 6, 14, 5, 11, 10\},$$

$$c_2 = \{21, 17, 22\}, c_3 = \{24\},$$

它们也是相关足码的集合：

$$c_1 = \{(2, 1), (3, 2), (4, 3), (3, 1), (4, 2), (4, 1)\},$$

$$c_2 = \{(6, 4), (5, 4), (6, 5)\}, c_3 = \{(6, 5)\}.$$

我们先来讨论 c_1 中的点的归属问题。在 c_1 中任选一个关连号码设为 $(4, 3)$ 的点，如把它放入 N_4 层，就说这个点是按其第 1 足码决定其归属的；如放入 N_3 层，则是按第 2 足码决定归属的。因为 $(4, 3)$ 这个点必定是 N_3 或 N_4 的邻点，所以决不应把它放入 N_3, N_4 以外的层中去，如果 c_1 中的点全部按第 1 足码决定其归属，则原来的 N_1, N_2, N_3, N_4 将扩充成：

$$N_1 = \{1\}, N_2 = \{3, 4, 2\},$$

$$N_3 = \{7, 8, 9, 6, 5\}, N_4 = \{12, 13, 15, 16, 14, 11, 10\}.$$

按照前面提到的应使诸层最大宽度最小的原则，我们应按第 2 种方式来决定 c_1 中各点的归属，类似的讨论可施之于 c_2 和 c_3 。当 c_1, c_2, c_3 中各点都已分配完毕后，即可逐层进行点的编号 N_1, N_2, \dots, N_6 各层中的点的划分情况及重新编号情况如图 6.10 和 6.11 相应的图，这时的带宽已缩减到 $\beta = 5$ ，外形则已减到 98。

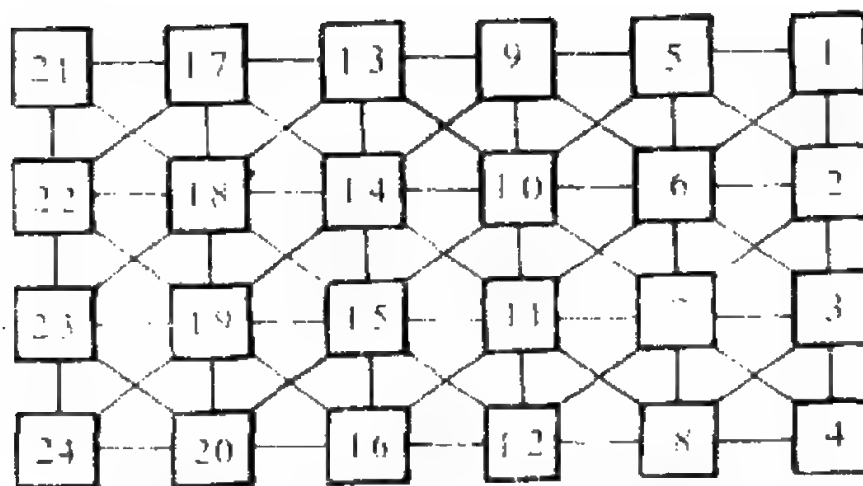


图6.11 按GPS方法决定顶点编号

§3 带状方程组解法

方程组的系数阵是带状阵的情形，在实用中较为普遍，且当联系到内存及计算量的节省时，有些实用性的技巧值得提出来，本节着重讨论这方面的问题。

3.1 等带宽矩阵的消去法

本段谈到的消去法是指列主元消去法，系数阵本身可以不必是对称的。

设

$$Ax = b \quad (6.3.1)$$

为所求方程组，系数阵 A 的带宽为 $2m+1$ 。当 A 以压缩形式存于二维数组 $L[1:n, 1:2m+1]$ 时，记

$$L[i', j'] = A[i, j],$$

$$i' = i, j' = j - i + m + 1, \quad (6.3.2)$$

$$i = 1, 2, \dots, n, \max(1, i - m) \leq j \leq \min(i + m, n)$$

当 $n = 6, m = 2$ 时，带形区域内的元素在矩阵 A ， L 中的存贮情况见图6.12。

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & & & \\ a_{21} & a_{22} & a_{23} & a_{24} & & \\ a_{31} & a_{32} & a_{33} & a_{34} & a_{35} & \\ & a_{42} & a_{43} & a_{44} & a_{45} & a_{46} \\ & & a_{53} & a_{54} & a_{55} & a_{56} \\ & & & a_{64} & a_{65} & a_{66} \end{pmatrix}$$

$$L = \begin{pmatrix} \times & \times & a_{11} & a_{12} & a_{13} \\ \times & a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} & a_{35} \\ a_{42} & a_{43} & a_{44} & a_{45} & a_{46} \\ a_{53} & a_{54} & a_{55} & a_{56} & \times \\ a_{64} & a_{65} & a_{66} & \times & \times \end{pmatrix}$$

图6.12 带状阵及其压缩存储

图中画“×”处的元素可以是任意数。由于消元将在 L 区域上进行，但原矩阵中的第1列的元素在 L 中是被排在第 $m+1$ 条上行对角线上的，因此在第一次消元之前，应先在前 m 行中，分别将第1, 2, ..., m 行左移 $m, m-1, \dots, 1$ 次，使 L 的前 m 行还原成矩阵原来的形状，其后 m 行中原来画“×”处的元素亦应同时清零，即使 L 在第一次消元以前的初态为

$$L = \begin{pmatrix} a_{11} & a_{12} & a_{13} & 0 & 0 \\ a_{21} & a_{22} & a_{23} & a_{24} & 0 \\ a_{31} & a_{32} & a_{33} & a_{34} & a_{35} \\ a_{42} & a_{43} & a_{44} & a_{45} & a_{46} \\ a_{53} & a_{54} & a_{55} & a_{56} & 0 \\ a_{64} & a_{65} & a_{66} & 0 & 0 \end{pmatrix} \quad \left. \begin{matrix} \\ \\ \\ \\ \end{matrix} \right\} (m+1)$$

当在 L 的前 $m+1$ 行的第1列中进行列主元消去法时，主

元的选取、换行以及实现消元的步骤与通常的列主元消去法步骤相同。第1次消元结束时， L 的前 $m+1$ 行的形状及数值意义为

$$\begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1,m+1}^{(1)} \\ 0 & a_{22}^{(1)} & \cdots & a_{2,m+1}^{(1)} \\ \cdots & \cdots & \cdots & \cdots \\ 0 & a_{m+1,1}^{(1)} & \cdots & a_{m+1,m+1}^{(1)} \end{pmatrix},$$

此时第1行中的信息即第一次消元结束时的输出，它们在回代时还要用到，应予保留，而当把第2~ $m-1$ 行左移一列时， L 在第2~ $m+2$ 行上的元素正好合于第2次消元时的需要，之后，消元即应在第2~ $m+2$ 行上进行，仿此直到消元过程进行完毕。当然，在行与行间进行换行及消元步骤时，对方程组的右端项向量 b 亦应进行相应的处理，第 i 次消元时，参加消元的行数一般应是 $m+1$ 行，即第 i 行到第 $i+m$ 行，当 $i+m > n$ 时，则应是第 i 行到第 n 行，参加消元的列数可以一直是 $2m+1$ 列，消元结束时，开始回代，而形成的解在 b 中，算法的具体步骤如下：

算法6.3.1 等带宽阵的列主元消去法：

1) 置 $m_1 = m + 1$, $m_2 = 2m + 1$;

2) 对 $r = 1, 2, \cdots, m_1 - 1$,

2.1) 对 $i = 1, \cdots, m_1 - r$,

2.1.1) 置 $L[r, j-1] = L[r, j] (j = 1, \cdots, m_2)$,

2.1.2) 置 $L[r, m_2] = L[n+1-r, m_2-i+1]$
 $= 0$,

2.1.3) NEXT i ,

2.2) NEXT r ,

3) 对 $i = 1, \cdots, n-1$,

- 3.1) 置 $k = i$,
- 3.2) 对 $r = i + 1, \dots, m_1$,
- 3.3) 如果 $|L[r, 1]| > |L[k, 1]|$ 则置 $k = r$,
- 3.4) NEXT r ;
- 4) 如果 $k \neq i$ 则转 5), 否则转 7);
- 5) $b[i] \leftrightarrow b[k]$;
- 6) $L[i, j] \leftrightarrow L[k, j] \quad j = 1, \dots, m_2$;
- 7) 置 $b[i] = b[i]/L[i, 1]$;
置 $L[i, j] = L[i, j]/L[i, 1], (j = 2, \dots, m_2)$ 。
- 8) 对 $r = i + 1, \dots, m_1$,
- 8.1) 置 $t = L[r, 1], b[r] = b[r] - t * b[i]$;
- 8.2) $L[r, j - 1] = L[r, j] - t * L[i, j], (j = 2, \dots, m_2)$;
- 8.3) 置 $L[r, m_2] = 0$;
- 8.4) NEXT r ;
- 9) 如果 $m_1 \neq n$ 则置 $m_1 \leftarrow m_1 + 1$;
- 10) NEXT i ;
- 11) 置 $b[n] = b[n]/L[n, 1], j_m = 2$;
- 12) 对 $r = n - 1, n - 2, \dots, 1$,
- 12.1) $b[r] = b[r] - L[r, j_m] * b[r - 1 + j_m]$,
($j = 2, \dots, j_m$),
- 12.2) 如果 $j_m \neq m$ 则置 $j_m \leftarrow j_m + 1$;
- 12.3) NEXT r 。

3.2 对称等带宽矩阵的消去法

由于不带行交换的高斯消去法在整个消元过程中, 每个行上第一个非零元左侧的元素将始终保持是零, 因而当矩阵 A 是正定对称等带宽矩阵时, 则在分解式

$$A = LDL^T \quad (6.3.3)$$

中， L 将与 A 有同样的带宽，这里 L 是下三角阵， D 是对角阵，它们都是非奇异的。当规定

$$d_{i,i} = 1/l_{i,i} \quad (6.3.4)$$

时，分解式 (6.3.3) 是唯一的。

设 m 为 A 的半带宽，比较 (6.3.3) 两端的下三角部分，有

$$\begin{aligned} a_{ij} &= \sum_{k=1}^n l_{ik} d_{kk} l_{jk} \quad (j \leq i) \\ &= \sum_{k=1}^j l_{ik} l_{jk} / l_{kk}, \end{aligned}$$

得

$$\begin{aligned} l_{i,i} &= a_{i,i} - \sum_{k=1}^{i-1} l_{ik}^2 / l_{kk} \\ l_{r,i} &= a_{r,i} - \sum_{k=1}^{i-1} l_{rk} l_{ik} / l_{kk} \\ t_r &= \max(1, r-m). \end{aligned} \quad (6.3.5)$$

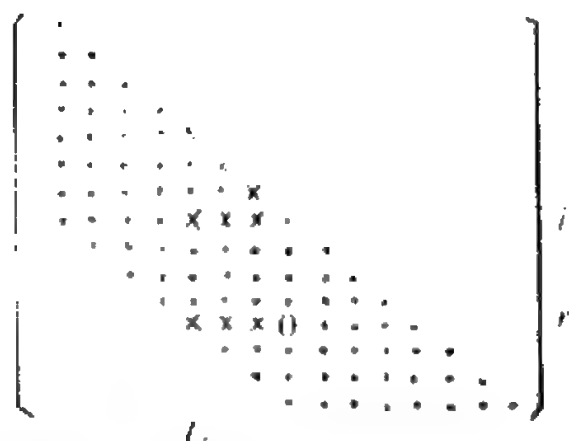


图6.13 计算 $l_{r,i}$ 所需的信息

由图 6.13 中所示的计算 $l_{r,i}$ 所需的信息，包括画“ \times ”处的元素，这是 L 中在计算 $l_{r,i}$ 之前已经算出的信息；画“ \circ ”处的元素是 A 阵中的 $a_{r,i}$ ，可见 L 中的元素可以逐列按由上而下的次序也可以逐行按自左至右的次序进行计算。当新的 $l_{r,i}$ 算出以后，即可用它来取代 $a_{r,i}$ 的存贮位置。由

$$LDL^T x = b, L^T x = \bar{b}, LD \bar{b} = b,$$

可得

$$\begin{aligned}\bar{b}_i &= b_i - \sum_{k=m_i}^{i-1} l_{ik} \bar{b}_k / l_{kk} \\ x_i &= \left(\bar{b}_i - \sum_{k=i}^{p_i-1} l_{ki} x_k \right) / l_{ii}\end{aligned}\quad (6.3.6)$$

$$m_i = \max(1, i - m), \quad p_i = \min(n, i + m).$$

当 A 的下三角部分以压缩形式存于二维数组 $L[1:n, 1:m+1]$ 时, 其地址检索关系为:

$$\begin{aligned}A[i, j] &= L[i', j'], \\ i' &= i, j' = j - i + m + 1.\end{aligned}\quad (6.3.7)$$

于 (6.3.5)、(6.3.6) 中换 a_{ri}, l_{ik}, \dots 为 $L[r, i - r + m + 1], L[i, k - i + m + 1], \dots$, 则得压缩存贮形式下的矩阵分解及求解的计算公式。

算法 6.3.2 对称等带宽矩阵的消去法:

- 1) 置 $m_1 = m + 1$;
- 2) 对 $i = 1, 2, \dots, n$,
 - 2.1) 置 $t = \max(1, i - m)$,
 - 2.2) 对 $j = t, t + 1, \dots, i$,
 - 2.2.1) $L[i, j - i + m_1] \leftarrow L[i, j - i + m_1]$

$$- \sum_{k=t}^{j-1} L[i, k - i + m_1] L[j, k - j + m_1] / L[k, m_1],$$
 - 2.3) $b[i] \leftarrow \left(b_i - \sum_{k=i}^{t-1} L[i, k - i + m_1] b[k] \right) / L[i, m_1];$
- 3) NEXT i ;
- 4) 对 $i = n, n - 1, \dots, 1$,
 - 4.1) 置 $t = \min(n, i + m)$,

$$4.2) \quad b[i] = \left(b[i] - \sum_{k=i+1}^i L[k, i-k+m_1] * b[k] \right) / L[i, m_1];$$

5) NEXT i.

3.3 对称变带宽矩阵的消去法

当 A 是对称正定的变带宽矩阵时, 其分解式

$$A = LDL^T$$

中, L 在每个行上的半带宽应与 A 在同行上的一致, 这里 L 、 D 的意义与 3.2 段中的规定是一样的.

在 A 是变带宽的情形下, L 的计算公式与 A 是等带宽情形下的区别仅在于: 应把公式 (6.3.5) 中求和符号的下限 $t_r = \max(1, r-m)$ 换为

$$t_r = \max(m_i, m_r) \quad (i \leq r), \quad (6.3.8)$$

这里 m_i, m_r 分别是第 i 行第 r 行上第 1 个非零元的列号, 它们可分别用第 i, r 行上的半带宽计算出来. 于是, 可写出相应于 (6.3.5) 的公式:

$$\begin{aligned} l_{ii} &= a_{ii} - \sum_{k=t_i}^{i-1} l_{ik}^2 / l_{kk}, \\ l_{ri} &= a_{ri} - \sum_{k=t_r}^{i-1} l_{rk} l_{ik} / l_{kk}, \\ t_r &= \max(m_i, m_r), \\ m_i &= i - \beta_i, m_r = r - \beta_r, \end{aligned} \quad (6.3.9)$$

和

$$\begin{aligned} \bar{b}_i &= b_i - \sum_{k=m_r}^{i-1} l_{ik} \bar{b}_k / l_{kk}, \\ x_i &= \left(\bar{b}_i - \sum_{k=i}^n l_{ki} x_k \right) / l_{ii}. \end{aligned} \quad (6.3.10)$$

当变带宽矩阵的下三角部分以压缩形式存放在一维数组 $L[1:S]$ 时, 据地址检索关系:

$$\begin{aligned} L(s) &= A[i, j], \\ s &= d[i] + j - i, \end{aligned} \quad (6.3.11)$$

$$i = 1, 2, \dots, n, j = m_i, \dots, i.$$

这里

$$s = \sum_{i=1}^n (\beta_i + 1), \quad (6.3.12)$$

而 $d(k)$ 是 A 的第 k 个对角元在 L 数组中的位置. 对于 (6.3.9)、(6.3.10), 分别将 l_{ik}, l_{kk}, \dots 换为 $L[d[i] + k - i], L[d[k]]$, 即可得出对称变带宽存贮下的矩阵分解及方程组求解的计算公式. 注意到公式 (6.3.10) 中计算 x_i 时要用到 L 在第 i 列中的数据, 而 L 区域的形状对列是非凸的, 因此检索关系 (6.3.11) 只有当 (i, j) 属于变带宽区域时才能使用, 否则 l_{ij} 便应是零. 判别 (i, j) 是否属于变带宽区域只需看如下不等式

$$m_i = i - d[i] + d[i - 1] + 1 \leq j \leq i \quad (6.3.13)$$

是否成立即可.

算法 6.3.3 变带宽对称矩阵的消去法:

1) 对 $i = 1, 2, \dots, n$,

1.1) 置 $i_0 = d[i] - i, m_i = i - (d[i] - d[i - 1]) + 1$,

1.2) 对 $j = m_i, \dots, i$,

1.2.1) 置 $j_0 = d[j] - j, m_j = j - (d[j] - d[j - 1]) + 1, m_{ij} = \max(m_i, m_j), ij = i_0 + j$,

1.2.2) $L[ij] \leftarrow L[ij] - \sum_{k=m_{ij}}^{j-1} L[i_0 + k] L[j_0 + k] / L[d[k]],$

1.2.3) NEXT j ,

1.3) $b[i] \leftarrow (b[i] - \sum_{k=m_i}^{i-1} L[i_0 + k] b[k]) / L[d[k]],$

1.4) NEXT i ,

2) $b[n] \leftarrow b[n] / L[d[n]],$

3) 对 $i = n - 1, \dots, 1$,

3.1) 置 $s = 0$,

- 3.2) 对 $k = i + 1, \dots, n$,
- 3.2.1) 置 $k_0 = d[k] - k$,
 $m_k = k - (d[k] - d[k - 1]) + 1$,
- 3.2.2) 如果 $m_k \leq i \leq k$, 则
 $s \leftarrow s + L[k_0 + i]b[k]$,
- 3.2.3) NEXT k ,
- 3.3) $b[i] \leftarrow (b[i] - S) / L[d[i]]$,
- 4) NEXT i .

第六章 习 题

- 6.1 当按压缩存储方式把对称正定矩阵存入内存时, 试
- 1) 建立压缩存储矩阵与原矩阵的地址检索关系,
 - 2) 导出用压缩存储矩阵元素表示的有关消元及回代过程公式.
- 6.2 证明, 在对称正定带状阵的 Cholesky 分解

$$A = LL^T$$

中, L 的任意一行上的半带宽与 A 的同一行上的半带宽相同.

6.3 在链式存储中, 如果因行列信息出错不能检索出元素的地址, 试编出能诊断这一情况并予以警告的程序.

6.4 在链式存储中, 设欲在第 j 列增加的元素是本列的第一个非零元或最末一个非零元, 试写出相应的检索步骤.

6.5 在链式存储中, 设第 j 列的某元素的行号为 i , 应依怎样的步骤找出同列中离它最近的非零元信息条首址?

6.6 设图 6.14 为某一对称正定矩阵的无向图为

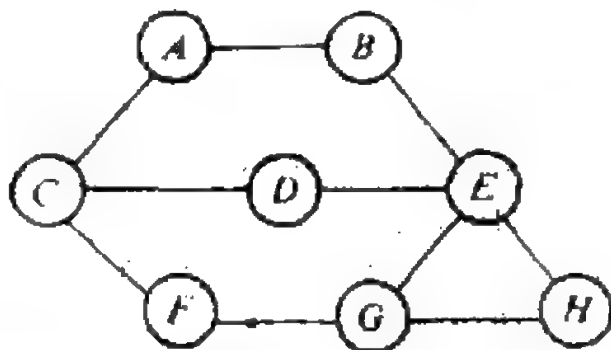


图 6.14

试分别以 $x_C, x_D, x_G, x_B, x_H, x_E, x_A, x_F$ 及 $x_H, x_F, x_D, x_C, x_E, x_B, x_A$ 为次序写出相应矩阵的 Ciolesky 分解的下三角阵中的非零元分布的情况, 两种排序中以何者为好? 能否据此提出一般排序的建议?

6.7 用 Tewarsen 方法编制第 k 次消元过程中的求列主元的子程序。

6.8 用 Markowitz 方法编制第 k 次消元过程中的求列主元的子程序。

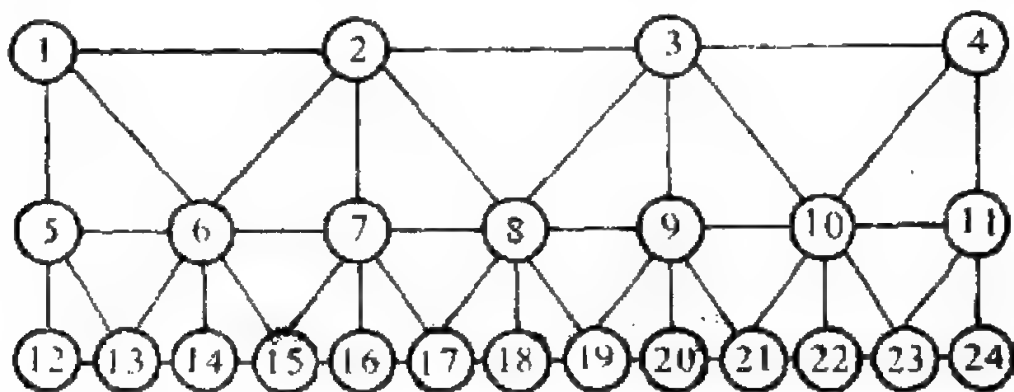


图6.15

6.9 按图右中节点编号的次序写出相应的对称正定矩阵各行上的半带宽并求其外形。

6.10 据上图以节点 12 为根写出分层结构中各层的节点, 按 CM 方法和 RCM 方法重新排序, 求出两种排序下的矩阵的带宽和外形。

参 考 书

1. R. P. Tewarson, "Sparse Matrices", Academic Press, New York (1973).
2. J. K. Reid, "Sparse Matrices" 引自 "The State of the art in numerical Analysis", Academic Press, London, New York San Francisco, (1977).
3. A. Jennings, "Matrix Computation for Engineers and Scientists", John Wiley & Sons, London, New York, Sydney, Toronto (1977).
4. Iain S. Duff "稀疏矩阵研究的综述", 应用数学与计算数学 1980 年第 1, 2 期; 唐路新译自 Proc. IEEE, 1977, 65, No. 4.
5. "数值代数的一些发展概况" 中国科学院计算中心图书情报室 (1978).

第七章 线性方程组的迭代解法

§1 引言

线性方程组的直接解法, 对于低阶方程组是很有成效的。对于高阶稀疏方程组, 目前虽然也介绍了各种有效的方法, 但是由于有存贮量大, 编程序的技巧要求较高等缺陷, 在实践中不免会遇到一些困难。迭代法则是克服这方面缺点的有效方法。迭代法与直接法的根本区别在于, 后者即使每一步的计算都是精确的运算, 得到的仍是近似解, 而前者则是从一系列的近似解中去获得满足精度要求的近似解的。由于迭代法必须考虑收敛性和误差, 所以需要根据问题的背景和条件去选择合适的方法, 即收敛速度较快的方法。这是研究迭代法的两个重要问题。

现在我们以三阶方程组为例, 对数值分析中介绍过的迭代方法进行回顾。

设三阶方程组为

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 = b_1, \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 = b_2, \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 = b_3. \end{cases} \quad (7.1.1)$$

设 $a_{ii} \neq 0$ ($i = 1, 2, 3$), 将 (7.1.1) 变形为

$$\begin{cases} x_1 = -\frac{a_{12}}{a_{11}}x_2 - \frac{a_{13}}{a_{11}}x_3 + \frac{b_1}{a_{11}}, \\ x_2 = -\frac{a_{21}}{a_{22}}x_1 - \frac{a_{23}}{a_{22}}x_3 + \frac{b_2}{a_{22}}, \\ x_3 = -\frac{a_{31}}{a_{33}}x_1 - \frac{a_{32}}{a_{33}}x_2 + \frac{b_3}{a_{33}}. \end{cases} \quad (7.1.2)$$

把它写成矩阵形式，则为

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} a_{11}^{-1} & & \\ & a_{22}^{-1} & \\ & & a_{33}^{-1} \end{pmatrix} \begin{pmatrix} 0 & -a_{12} & -a_{13} \\ -a_{21} & 0 & -a_{23} \\ -a_{31} & -a_{32} & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} a_{11}^{-1} & & \\ & a_{22}^{-1} & \\ & & a_{33}^{-1} \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}.$$

将(7.1.2)写成迭代格式，就得到简单迭代法的迭代格式：

$$\begin{cases} x_1^{(m)} = -\frac{a_{12}}{a_{11}}x_2^{(m-1)} - \frac{a_{13}}{a_{11}}x_3^{(m-1)} + \frac{b_1}{a_{11}}, \\ x_2^{(m)} = -\frac{a_{21}}{a_{22}}x_1^{(m-1)} - \frac{a_{23}}{a_{22}}x_3^{(m-1)} + \frac{b_2}{a_{22}}, \\ x_3^{(m)} = -\frac{a_{31}}{a_{33}}x_1^{(m-1)} - \frac{a_{32}}{a_{33}}x_2^{(m-1)} + \frac{b_3}{a_{33}}. \end{cases} \quad (7.1.3)$$

对于迭代格式(7.1.3)，还可以进行改进。由于在做(7.1.3)的第2式时， $x_1^{(m)}$ 已经计算出来了，在做第3式时， $x_2^{(m)}$ 也计算出来了，因此，我们可以将(7.1.3)改写成

$$\begin{cases} x_1^{(m)} = -\frac{a_{12}}{a_{11}}x_2^{(m-1)} - \frac{a_{13}}{a_{11}}x_3^{(m-1)} + \frac{b_1}{a_{11}}, \\ x_2^{(m)} = -\frac{a_{21}}{a_{22}}x_1^{(m)} - \frac{a_{23}}{a_{22}}x_3^{(m-1)} + \frac{b_2}{a_{22}}, \\ x_3^{(m)} = -\frac{a_{31}}{a_{33}}x_1^{(m)} - \frac{a_{32}}{a_{33}}x_2^{(m)} + \frac{b_3}{a_{33}}. \end{cases} \quad (7.1.4)$$

这就是GS法的迭代格式。如果在这个迭代格式中引进一个加速迭代参数，即所谓松弛因子，并将(7.1.4)改写成

$$\begin{cases} x_1^{(m)} = \omega \left(-\frac{a_{12}}{a_{11}} x_2^{(m-1)} - \frac{a_{13}}{a_{11}} x_3^{(m-1)} + \frac{b_1}{a_{11}} \right) + (1-\omega) x_1^{(m-1)}, \\ x_2^{(m)} = \omega \left(-\frac{a_{21}}{a_{22}} x_1^{(m)} - \frac{a_{23}}{a_{22}} x_3^{(m-1)} + \frac{b_2}{a_{22}} \right) + (1-\omega) x_2^{(m-1)}, \\ x_3^{(m)} = \omega \left(-\frac{a_{31}}{a_{33}} x_1^{(m)} - \frac{a_{32}}{a_{33}} x_2^{(m)} + \frac{b_3}{a_{33}} \right) + (1-\omega) x_3^{(m-1)}. \end{cases} \quad (7.1.5)$$

这就是逐次松弛法的迭代格式。可简记成SOR法。如果设

$$D = \begin{pmatrix} a_{11} & & \\ & a_{22} & \\ & & a_{33} \end{pmatrix}, \quad L = \begin{pmatrix} 0 & 0 & 0 \\ -a_{21} & 0 & 0 \\ -a_{31} & -a_{32} & 0 \end{pmatrix},$$

$$U = \begin{pmatrix} 0 & -a_{12} & -a_{13} \\ 0 & 0 & -a_{23} \\ 0 & 0 & 0 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix},$$

则简单迭代法的迭代格式为

$$\mathbf{x}^{(m)} = D^{-1}(L + U)\mathbf{x}^{(m-1)} + D^{-1}\mathbf{b}, \quad (7.1.6)$$

GS法的迭代格式为

$$\mathbf{x}^{(m)} = D^{-1}\{L\mathbf{x}^{(m)} + U\mathbf{x}^{(m-1)}\} + D^{-1}\mathbf{b},$$

或 $(I - D^{-1}L)\mathbf{x}^{(m)} = D^{-1}U\mathbf{x}^{(m-1)} + D^{-1}\mathbf{b}$

即 $\mathbf{x}^{(m)} = (I - D^{-1}L)^{-1}D^{-1}U\mathbf{x}^{(m-1)} + (I - D^{-1}L)^{-1}D^{-1}\mathbf{b}, \quad (7.1.7)$

SOR法的迭代格式为

$$\begin{aligned} \mathbf{x}^{(m)} = & \omega(D^{-1}L\mathbf{x}^{(m)} + D^{-1}U\mathbf{x}^{(m-1)} \\ & + D^{-1}\mathbf{b}) + (1-\omega)\mathbf{x}^{(m-1)}, \end{aligned}$$

或 $(I - \omega D^{-1}L)\mathbf{x}^{(m)} = \omega D^{-1}U\mathbf{x}^{(m-1)} + (1-\omega)\mathbf{x}^{(m-1)} + \omega D^{-1}\mathbf{b}$

即 $\mathbf{x}^{(m)} = (I - \omega D^{-1}L)^{-1}\{(1-\omega)I + \omega D^{-1}U\}\mathbf{x}^{(m-1)} + (I - \omega D^{-1}L)^{-1}\omega D^{-1}\mathbf{b}. \quad (7.1.8)$

当 $\omega=1$ 时, SOR法化为GS法。

可以把三种迭代格式写成统一的形式:

$$\mathbf{x}^{(m)} = B\mathbf{x}^{(m-1)} + \mathbf{g}, \quad (7.1.9)$$

其中矩阵 B 称为迭代法 (7.1.9) 中的**迭代矩阵**.

对于 n 阶方程组

$$A\mathbf{x} = \mathbf{b}, \quad (7.1.10)$$

其中 A 是 n 阶方阵, \mathbf{x} , \mathbf{b} 为 n 维列向量; 如果令 $A = D - L - U$,

其中

$$D = \begin{pmatrix} a_{11} & & & & \\ & & & 0 & \\ & & a_{22} & & \\ & 0 & & \ddots & \\ & & & & a_{nn} \end{pmatrix},$$

$$L = \begin{pmatrix} 0 & & & & \\ & & & & \\ & -a_{21} & & & \\ \vdots & & \ddots & & \\ -a_{n1} & \cdots & -a_{n,n-1} & & 0 \end{pmatrix},$$

$$U = \begin{pmatrix} & & & & \\ & & & & \\ & -a_{12} & \cdots & -a_{1n} & \\ & 0 & \ddots & \vdots & \\ & 0 & \ddots & -a_{n-1,n} & \\ & & & 0 & \end{pmatrix},$$

则迭代法 (7.1.6)、(7.1.7)、(7.1.8)、(7.1.9) 同样成立.

对于任意的初始值 $\mathbf{x}^{(0)}$, 由 (7.1.9) 得出的向量序列 $\mathbf{x}^{(1)}$, $\mathbf{x}^{(2)}$, ..., $\mathbf{x}^{(n)}$, ..., 如果收敛于某个向量 \mathbf{x}^* , 则 \mathbf{x}^* 就是方程组 (7.1.1) 的精确解. 因此, 迭代法实际上是一种从已有近似解去计算新近似解的法则. 从 (7.1.9) 看出, 新近似解 $\mathbf{x}^{(m)}$ 是已

有近似解 $x^{(m-1)}$ 的线性函数，这种迭代法叫**线性迭代法**。在本章的最后一节我们还将介绍**非线性迭代法**。此外，由于(7.1.9)中 $x^{(m)}$ 仅是前一次近似解 $x^{(m-1)}$ 的函数，这种迭代法又叫做**一阶迭代法**。如果不仅依赖前一次的近似解，而且还依赖前面相邻 K 次的近似解 $x^{(m-1)}, x^{(m-2)}, \dots, x^{(m-K)}$ ，则称为 **K 阶迭代法**。例如 Chebyshev 半迭代法就是二阶迭代法。从(7.1.9)中还可看出，无论 m 为何值，计算 $x^{(m)}$ 的规则均是不变的，这种迭代法称为**定常迭代法**。Chebyshev 半迭代法还是一种非定常的线性迭代法。

我们构造迭代格式的目的是要求原方程组的解。因此，在构造迭代格式时，必须注意迭代法的极限是否是原方程组的解。如果对于任取的初始向量 $x^{(0)}$ ，由(7.1.9)所产生的向量序列 $\{x^{(m)}\}$ 都有相同的极限，并且其极限就是原方程组的解 x^* ，则称迭代格式(7.1.9)是**收敛的**，否则称为**不收敛的**。

定理7.1.1 对于任意右端向量 g 和初始向量 $x^{(0)}$ ，迭代法(7.1.9)收敛的充要条件是 $S(B) < 1$ 。这里 $S(B)$ 表示 B 的谱半径。

证明 必要性。设迭代法(7.1.9)收敛且极限是 x^* ，即

$$x^* = Bx^* + g.$$

将上式与(7.1.9)相减得到

$$x^{(m)} - x^* = B(x^{(m-1)} - x^*).$$

令

$$r^{(m)} = x^{(m)} - x^*, \quad (7.1.11)$$

则有

$$r^{(m)} = Br^{(m-1)} = B^2 r^{(m-2)} = \dots = B^m r^{(0)}. \quad (7.1.12)$$

据假设

$$\lim_{m \rightarrow \infty} r^{(m)} = \lim_{m \rightarrow \infty} B^m r^{(0)} = 0,$$

故必有 $\lim_{m \rightarrow \infty} B^m = 0$ ，因此 $S(B) < 1$ 。

充分性。设 $S(B) < 1$, 则 $I - B$ 非奇异, 从而方程组

$$(I - B)x = g$$

有唯一解, 记为 x^* 。由 $S(B) < 1$, 可得 $\lim_{m \rightarrow \infty} B^m = 0$, 故有

$$\lim_{m \rightarrow \infty} r^{(m)} = \lim_{m \rightarrow \infty} B^m r^{(0)} = 0,$$

即 $\lim_{m \rightarrow \infty} x^{(m)} = x^*$ 。

对于迭代法, 不但要知道它是否收敛, 而且还必须研究它的收敛速度。

由 (7.1.12), 可得

$$\|r^{(m)}\| \leq \|B^m\| \|r^{(0)}\|,$$

因此, $\|B^m\|$ 的大小可以决定 $r^{(m)}$ 收敛于零向量的速度。如果要使 $\|r^{(m)}\|$ 比 $\|r^{(0)}\|$ 减小 e 倍 ($0 < e < 1$)

即 $\|r^{(m)}\| \leq e \|r^{(0)}\|$,

为此仅要求

$$\|B^m\| \leq e \text{ 或 } (\|B^m\|^{1/m})^m \leq e.$$

两边取对数, 得到

$$\ln e \geq \min (\|B^m\|^{1/m}),$$

即 $m \geq -\ln e / (-\frac{1}{m} \ln \|B^m\|)$. (7.1.13)

由此可见, 要达到精度为 e 的要求, 所需迭代的最少次数与 $(-\ln \|B^m\|/m)$ 成反比, 这个量愈大, 所需的迭代次数就愈小。因此, 我们可以用这个量来衡量迭代格式的收敛速度。

定义 7.1.1 称 $R_m(B) = -\frac{1}{m} \ln \|B^m\|$ 为迭代法 (7.1.9) 的**平均收敛速度**。

由于平均收敛速度难于求出, 因此我们往往用一粗略的估计式来代替它。若 B 仅是有线性初等因子的矩阵, 则必存在 B 的线性无关的特征向量系 x_1, x_2, \dots, x_n 和相应的特征值 $\lambda_1, \lambda_2,$

..., λ_n . 对于任意初始误差向量 $r^{(0)}$, 设

$$r^{(0)} = \sum_{i=1}^n a_i x_i,$$

代入 (7.1.12) 得到

$$\begin{aligned} r^{(m)} &= \sum_{i=1}^n a_i B^m x_i = \sum_{i=1}^n a_i \lambda_i^m x_i \\ &= (S(B))^m \sum_{i=1}^n \left(\frac{\lambda_i}{S(B)} \right)^m a_i x_i, \end{aligned}$$

两边取范数, 便有

$$\|r^{(m)}\| \leq (S(B))^m \left\| \sum_{i=1}^n \left(\frac{\lambda_i}{S(B)} \right)^m a_i x_i \right\|.$$

当 m 很大时, $\|r^{(m)}\|$ 收敛于零的速度取决于 $(S(B))^m$ 的大小, 如果要求 $\|r^{(m)}\| \leq \varepsilon$, 那么, 迭代的次数 m 必须近似地满足

$$(S(B))^m \leq \varepsilon,$$

$$\text{即} \quad m \geq (-\ln \varepsilon) / (-\ln(S(B))), \quad (7.1.14)$$

因而所需要的迭代次数 m 与 $(-\ln(S(B)))$ 成反比. 同时, 还可以证明 $\lim_{m \rightarrow \infty} (\|B^m\|)^{\frac{1}{m}} = S(B)$.

定义 7.1.2 称 $R(B) = -\ln(S(B))$ 为迭代法 (7.1.9) 的渐近收敛速度. 简称收敛速度.

利用 $R(B)$ 可以比较各种迭代法收敛的快慢程度. 显然, 谱半径 $S(B)$ 愈小, 则收敛速度将愈快. 但是迭代次数的比较可靠的估计还是 (7.1.13), 而 (7.1.14) 仅是粗略的估计式. 特别, 当 B 的按模最大特征值所对应的 Jordan 块的阶数大于 1 时, 实际迭代的次数往往要比估计的次数大得多. 应用 (7.1.14) 时必须注意这一点 [参看 Iterative solution of large linear systems P85—P87.]

§2 逐次松弛法

逐次松弛法是一阶线性定常迭代法, 对于一类特殊的矩阵, 如果松弛因子 ω 选得恰当, 它的收敛速度会比 **Jacobi** 方法和 **GS** 方法快得多. 通常我们把 $\omega > 1$ 的松弛迭代法称为**超松弛迭代法**, 而把 $0 < \omega < 1$ 的松弛迭代法称为**低松弛迭代法**. 不加区别, 则统称为**逐次松弛法**. 松弛因子选在什么范围内松弛法才能收敛? 是我们最关心的问题之一. 下面我们将在数值分析的基础上, 较详细地来讨论逐次松弛法的收敛性和最佳松弛因子的选取问题.

2.1 逐次松弛法的收敛性

定理 7.2.1 当 ω 取实数时, 逐次松弛法 (7.1.8) 收敛的必要条件是 $0 < \omega < 2$. (7.2.1)

证明 设 (7.1.8) 的迭代矩阵为

$$J_\omega = (I - \omega D^{-1}L)^{-1}[(1 - \omega)I + \omega D^{-1}U], \quad (7.2.2)$$

则

$$\begin{aligned} \det(J_\omega) &= \det((I - \omega D^{-1}L)^{-1}) \det((1 - \omega)I + \omega D^{-1}U) \\ &= (1 - \omega)^n. \end{aligned}$$

因此, 矩阵 J_ω 的所有特征值之积为 $(1 - \omega)^n$, 故

$$S(J_\omega) \geq |1 - \omega|. \quad (7.2.3)$$

据定理 7.1.1. **SOR** 法收敛的充要条件为 $S(J_\omega) < 1$ 知, 必须 $|1 - \omega| < 1$, 即 $0 < \omega < 2$.

实用中碰到的系数矩阵, 多数是对称正定的. 例如, 用有限元素法求解弹性结构的静力平衡问题, 就可归结为解一个系数矩阵为对称正定的线性代数方程组的问题. 对于这类系数矩阵, 收敛性问题现已得到了较好的解决.

定理7.2.2 设 A 是实对称正定矩阵, 则 $0 < \omega < 2$ 时, 逐次松弛法是收敛的.

证明 因为 A 是实对称正定阵, 所以有

$$A = D - L - L^T, \quad (7.2.4)$$

$$J_\omega = (I - \omega D^{-1}L)^{-1}[(1 - \omega)I + \omega D^{-1}L^T]. \quad (7.2.5)$$

设 λ 为 J_ω 的任一特征值, \mathbf{x} 为相应的特征向量, 则

$$J_\omega \mathbf{x} = \lambda \mathbf{x},$$

即
$$[(1 - \omega)I + \omega D^{-1}L^T]\mathbf{x} = \lambda(I - \omega D^{-1}L)\mathbf{x},$$

$$(1 - \omega)\mathbf{x} + \omega D^{-1}L^T\mathbf{x} = \lambda(\mathbf{x} - \omega D^{-1}L\mathbf{x}).$$

两边用 $\mathbf{x}^H D$ 左乘, 得

$$(1 - \omega)\mathbf{x}^H D \mathbf{x} + \omega \mathbf{x}^H L^T \mathbf{x} = \lambda(\mathbf{x}^H D \mathbf{x} - \omega \mathbf{x}^H L \mathbf{x}). \quad (7.2.6)$$

由于 D 为对角矩阵且对角线元素均为正实数, 因此若令 $p = \mathbf{x}^H D \mathbf{x}$, 则 p 为正实数, 又 $(\mathbf{x}^H L^T \mathbf{x})^H = \mathbf{x}^H L \mathbf{x}$, 所以若令 $\mathbf{x}^H L \mathbf{x} = \alpha + i\beta$, 则 $\mathbf{x}^H L^T \mathbf{x} = \alpha - i\beta$. 代入 (7.2.6), 得到

$$(1 - \omega)p + \omega(\alpha - i\beta) = \lambda(p - \omega(\alpha + i\beta)).$$

因此

$$\lambda = \frac{(1 - \omega)p + \omega\alpha - i\omega\beta}{p - \omega\alpha - i\omega\beta},$$

$$|\lambda| = \frac{[p - \omega(p - \alpha)]^2 + \omega^2\beta^2}{(p - \omega\alpha)^2 + \omega^2\beta^2},$$

由于

$$(p - \omega(p - \alpha))^2 - (p - \omega\alpha)^2 = p\omega(2 - \omega)(2\alpha - p),$$

利用 A 的正定性可得

$$\mathbf{x}^H A \mathbf{x} = \mathbf{x}^H D \mathbf{x} - \mathbf{x}^H L \mathbf{x} - \mathbf{x}^H L^T \mathbf{x} = p - 2\alpha > 0.$$

从而, 当 $0 < \omega < 2$ 时, 便有

$$(p - \omega(p - \alpha))^2 - (p - \omega\alpha)^2 < 0,$$

得

$$|\lambda|^2 = \frac{(p - \omega(p - \alpha))^2 + \omega^2\beta^2}{(p - \omega\alpha)^2 + \omega^2\beta^2} < 1,$$

所以

$$S(J_\omega) < 1.$$

据定理 7.1.1. SOR 法收敛.

定义 7.2.1 设 A 为 n 阶方阵, 如对所有的 $i (1 \leq i \leq n)$ 均有

$$|a_{ii}| \geq \sum_{j=1, j \neq i}^n |a_{ij}| \quad (7.2.7)$$

成立, 则称 A 为**对角占优阵**. 如果至少有一个 i 值, (7.2.7) 有严格的不等号成立. 则称 A 为**弱角对角占优阵**. 若对所有的 i 值, (7.2.7) 均有严格不等号成立, 则称 A 为**严格对角占优阵**.

定义 7.2.2 如果存在 n 阶排列矩阵 P , 能使 n 阶方阵 A 化为:

$$PAP^T = \begin{pmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{pmatrix}, \quad (7.2.8)$$

其中 A_{11} 、 A_{22} 均为方阵, 则称 A 为**可约矩阵**, 否则称为**不可约矩阵**.

例如, 对于矩阵

$$A = \begin{pmatrix} 4 & 0 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & 0 & 2 & 0 \\ -1 & -1 & -1 & 1 \end{pmatrix},$$

存在矩阵

$$P = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix},$$

使得

$$PAP^T = \begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 3 & -1 & -1 \\ 0 & 0 & 2 & -1 \\ 0 & 0 & -1 & 4 \end{pmatrix},$$

故 A 为可约矩阵。

如果线性方程组 $Ax = b$ 的系数矩阵 A 为可约矩阵，那么，它可改写为

$$PAP^T(Px) = Pb, \quad (7.2.9)$$

若分别把 Px 和 Pb 记为 y 和 h ，并把 PAP^T 按 (7.2.8) 的形式进行分块，则 (7.2.9) 可以写成为

$$\begin{pmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} h_1 \\ h_2 \end{pmatrix}, \quad (7.2.10)$$

即系数矩阵为可约矩阵的方程组，可以约化为两个较低阶的方程组来求解。又如

$$A = \begin{pmatrix} 4 & -1 & 0 & 0 \\ -1 & 3 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

为不可约的弱对角占优阵，其 $\det(A) = 8 \neq 0$ 。不可约的弱对角占优阵是否一定为非奇异矩阵？可以证明如下定理（留作习题）。

定理 7.2.3 若 A 为不可约的弱对角占优阵，则 A 必为非奇异阵，即 $\det(A) \neq 0$ 。

定理 7.2.4 当系数矩阵 A 为不可约弱对角占优矩阵且松弛因子 ω 满足 $0 < \omega \leq 1$ ，则 SOR 法必定收敛。

证明 设 SOR 法的迭代矩阵为

$$J_\omega = (I - \omega D^{-1}L)^{-1}[(1 - \omega)I + \omega D^{-1}U], \quad (7.2.11)$$

其中 $A = D - L - U$ 须证 $S(J_\omega) < 1$ 。

用反证法，设 J_ω 有一个特征值 λ 满足 $|\lambda| \geq 1$ ，则 $J_\omega - \lambda I$ 应为奇异矩阵，即 $\det(J_\omega - \lambda I) = 0$ ，并且

$$\begin{aligned} J_\omega - \lambda I &= (I - \omega D^{-1}L)^{-1}[(1 - \omega)I + \omega D^{-1}U] - \lambda I \\ &= (I - \omega D^{-1}L)^{-1}[(1 - \omega)I + \omega D^{-1}U - \lambda(I - \omega D^{-1}L)] \end{aligned}$$

$$\begin{aligned}
&= (I - \omega D^{-1}L)^{-1}[(1 - \omega - \lambda)I + \omega D^{-1}U + \lambda\omega D^{-1}L] \\
&= (1 - \omega - \lambda)(I - \omega D^{-1}L)^{-1} \left[I - \frac{\omega}{\omega + \lambda - 1} D^{-1}U \right. \\
&\quad \left. - \frac{\lambda\omega}{\omega + \lambda - 1} D^{-1}L \right], \quad (7.2.12)
\end{aligned}$$

其中 $1 - \omega - \lambda \neq 0$ 和 $(I - \omega D^{-1}L)^{-1}$ 为非奇异矩阵. 由于 $A = D - L - U$ 为不可约弱对角占优阵, 所以, $D^{-1}A = I - D^{-1}L - D^{-1}U$ 也为不可约弱对角占优, 又由于矩阵 $D^{-1}A$ 中的零元素及其位置与矩阵 $I - \frac{\omega}{\omega + \lambda - 1} D^{-1}U - \frac{\lambda\omega}{\omega + \lambda - 1} D^{-1}L$ 中的零元素及其位置一样, 故后者也是不可约矩阵. 此外又因为

$$\begin{aligned}
\lambda &= (\lambda + \omega - 1) + (1 - \omega), \\
|\lambda| &\leq |\lambda + \omega - 1| + (1 - \omega), \\
|\lambda + \omega - 1| &\geq |\lambda| - (1 - \omega),
\end{aligned}$$

所以

$$\begin{aligned}
\left| \frac{\omega}{\lambda + \omega - 1} \right| &\leq \left| \frac{\lambda\omega}{\lambda + \omega - 1} \right| \leq \frac{|\lambda\omega|}{|\lambda| - (1 - \omega)} \\
&\leq \frac{|\lambda|\omega}{|\lambda| - (1 - \omega)|\lambda|} = 1,
\end{aligned}$$

于是
$$I - \frac{\omega}{\lambda + \omega - 1} D^{-1}U - \frac{\lambda\omega}{\lambda + \omega - 1} D^{-1}L \quad (7.2.13)$$

也是弱对角占优阵. 据定理 7.2.3, (7.2.13) 为非奇异矩阵, 故 (7.2.12) 的右端为非奇异矩阵. 这与左端为奇异矩阵相矛盾, 于是 $S(I_\omega) < 1$.

由于当 $\omega = 1$ 时, SOR 法即为 GS 法, 因此有下述推论.

推论 当系数矩阵 A 为不可约对角占优时, GS 方法收敛.

SOR 法的收敛问题与松弛因子 ω 的选择密切相关. 松弛因子选择得当, 可以大大提高收敛速度. 由于松弛法的应用比较广泛, 所以如何挑选最佳的松弛因子, 具有特别重要的意义.

2.2 最佳松弛因子的选取

前面我们已经证明，线性迭代法收敛的充要条件是，迭代矩阵的谱半径小于1，谱半径愈小，收敛就愈快。SOR法的迭代矩阵的谱半径为 $S(J_\omega)$ ，它是 ω 的函数。如何选取最佳松弛因子 ω_{opt} ，使得

$$S(J_{\omega_{opt}}) = \min_{\omega} S(J_{\omega})$$

是我们需要研究的课题。但是，对于一般的系数矩阵，目前尚无确定 ω_{opt} 的理论结果。实际计算时，大多是用计算经验来试算确定的。仅对某些特殊类型的矩阵，才有确定 ω_{opt} 的理论公式。在此，我们重点讨论具有所谓“性质A”的矩阵的最佳松弛因子的选取问题。最后，将介绍用试算法确定最佳松弛因子的方法。

定义7.2.3 若能将 n 个自然数的集合 $W = \{1, 2, \dots, n\}$ ，分成两个互不相交的子集 S_1 和 S_2 (即 $S_1 \cup S_2 = W$, $S_1 \cap S_2 = \emptyset$)，使得 n 阶矩阵 A 的所有非零非对角线元素 a_{ij} ($a_{ij} \neq 0$, $i \neq j$) 的足标对 (i, j) ，满足条件： $i \in S_1, j \in S_2$ 或 $j \in S_1, i \in S_2$ ，则称此矩阵具有“性质A”。

例如，将线性方程组

$$\begin{pmatrix} 1 & -0.30009 & 0 & -0.30898 \\ -0.30009 & 1 & -0.46691 & 0 \\ 0 & -0.46691 & 1 & -0.27471 \\ -0.30898 & 0 & -0.27471 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} 5.32088 \\ 6.07624 \\ -8.80455 \\ 2.67600 \end{pmatrix} \quad (7.2.14)$$

的系数矩阵记为 A ，取 $S_1 = \{1, 3\}$ ， $S_2 = \{2, 4\}$ ，容易验证矩阵 A 具有“性质A”。

如果把方程组(7.2.14)的第2, 3两行互相交换并把2, 3两列也互相交换, 则得

$$\begin{pmatrix} 1 & 0 & -0.30009 & -0.30898 \\ 0 & 1 & -0.46691 & -0.27471 \\ -0.30009 & -0.46691 & 1 & 0 \\ -0.30898 & -0.27471 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_3 \\ x_2 \\ x_4 \end{pmatrix} = \begin{pmatrix} 5.32008 \\ -8.80455 \\ 6.07624 \\ 2.67600 \end{pmatrix}.$$

若令 $S_1 = \{1, 2\}$, $S_2 = \{3, 4\}$, 仍可验证, 新系数矩阵也具有性质A.

定理 7.2.5 矩阵A具有“性质A”的充要条件是存在排列矩阵P, 使得

$$PAP^T = \begin{pmatrix} D_1 & R \\ S & D_2 \end{pmatrix}, \quad (7.2.15)$$

其中 D_1, D_2 是对角矩阵.

证明 必要性. 只需将具有“性质A”的矩阵A的号码属于 S_1 的各行排在前面, 把号码属于 S_2 的各行排在后面, 并把号码属于 S_1 的各列也排在前面, 把号码属于 S_2 的各列排在后面, 便得(7.2.15)型的矩阵.

充分性. 其证明作为习题.

对于线性方程组中的方程式和未知数, 在次序进行重排以后, 虽然方程组的精确解并不改变. 但是, 如果使用迭代法求解, 其收敛性和收敛速度却要引起变化. 同样对于具有性质A的矩阵, 次序经过相应的变化后, 虽然还保持“性质A”, 但其解的收敛性和收敛速度也要变化. 那末, 方程组和未知数的

次序应以怎样的排列为好呢？下面我们先看两个例子。

例7.2.1 对方程组(7.2.14)使用SOR法求解，取松弛因子 $\omega = 1.16$ ，实际迭代次数和结果见下表。

k	x_1	x_2	x_3	x_4	e_k
0	0	0	0	0	12.3095
1	6.1722	9.1970	-5.2320	3.6492	3.6659
2	9.6941	6.1177	-4.8999	4.4335	1.3313
3	8.3398	6.3188	-4.5914	3.9200	0.2302
4	8.4424	6.4480	-4.7151	4.0004	0.0767
5	8.5137	6.4202	-4.7068	4.0157	0.0288
6	8.4842	6.4253	-4.7005	4.0047	0.0051
7	8.4868	6.4288	-4.7031	4.0065	0.0016

例7.2.2 如果将方程组(7.2.14)的第3列与第4列交换，第3行与第4行交换，则得

$$\begin{pmatrix} 1 & -3.30009 & -3.30898 & 0 \\ -3.30009 & 1 & 0 & -0.46691 \\ -3.30898 & 0 & 1 & -0.27471 \\ 0 & -0.46691 & -0.27471 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_4 \\ x_3 \end{pmatrix} = \begin{pmatrix} 5.32088 \\ 6.07624 \\ 2.67600 \\ -8.80455 \end{pmatrix} \quad (7.2.16)$$

对这个方程组，用SOR方法求解，取松弛因子仍为 $\omega = 1.16$ ，其实际迭代次数及结果见下表。

K	x_1	x_2	x_3	x_4	e_i
0	0	0	0	0	12.3095
1	6.1722	-10.2133	3.6653	2.0618	6.2901
2	7.1996	-5.9369	5.7256	3.4629	1.9833
3	8.2639	-5.0442	6.2727	3.9047	0.4483
4	8.4330	-4.7645	6.3998	3.9837	0.0899
5	8.4785	-4.7153	6.4220	4.0031	0.0168
6	8.4859	-4.7050	6.4266	4.0059	0.0030
7	8.4874	-4.7032	6.4273	4.0065	0.0005

比较例7.2.1和例7.2.2不难看到, 例7.2.2中 **SOR** 法收敛速度较快, 这说明, 线性方程组的方程和未知数的次序的排列不同, **SOR**法的收敛速度也将不同, 那末, 例7.2.2的排列次序有什么特点呢? 为此, 还将研究一些有关概念.

定义7.2.4 若能将 n 个自然数的集合 $W = \{1, 2, \dots, n\}$ 分成 l 个互不相交的子集 S_1, S_2, \dots, S_l (即 $W = \bigcup_{k=1}^l S_k, S_i \cap S_j = \emptyset, i \neq j$), 使得矩阵 A 的所有非零的非对角元 $a_{ij} \neq 0 (i \neq j)$, 其足码对 (i, j) , 满足以下条件: 若 $i \in S_k$, 当 $j < i$ 时, $j \in S_{k-1}$, 而当 $j > i$ 时, $j \in S_{k+1}$, 则称此矩阵 A 具有**相容次序**.

若取 $S_1 = \{1\}, S_2 = \{2, 3\}, S_3 = \{4\}$, 容易验证(7.2.16)中的系数矩阵满足定义7.2.4, 即(7.2.16)中的系数矩阵是具有相容次序的. 方程组(7.2.16)的系数矩阵也具有“性质A”.

定理7.2.6 若矩阵具有相容次序, 则必具有“性质A”.

证明 若矩阵具有相容次序, 则存在满足定义7.2.4的子集 S_1, S_2, \dots, S_l , 令

$W_1 = \bigcup_{i \text{ 是奇数}} S_i$, $W_2 = \bigcup_{i \text{ 是偶数}} S_i$, 则 W_1 与 W_2 就是定义 7.2.3 的两个子集 S_1 和 S_2 , 即矩阵具有“性质 A”。

例 7.2.3 设矩阵 A 具有如下形式:

$$A = \begin{pmatrix} D_1 & U_1 & & \\ L_2 & D_2 & U_2 & \\ & \ddots & \ddots & \ddots \\ & & L_m & D_m \end{pmatrix} \quad (7.2.17)$$

其中 D_i ($i = 1, 2, \dots, m$) 为对角阵, L_{i+1}, U_i ($i = 1, 2, \dots, m-1$) 为相应阶数的长方阵。若将所有对角线子块 D_{2k} ($k = 1, 2, \dots$) 的行编号记为 S_1 , 将所有子块 D_{2k-1} ($k = 1, 2, \dots$) 的行编号记为 S_2 , 则容易验证, 集合 S_1 与 S_2 满足定义 7.2.3, 因此矩阵 A 具有“性质 A”。如果将 A 中属于对角线子块 D_k 的行编号记为 S_k , 显然, $\bigcup_{k=1}^m S_k = W$, $S_i \cap S_j = \emptyset$ ($i \neq j$), 而且容易验证这样划分 W 的子集是满足定义 7.2.4 的条件的, 所以矩阵 (7.2.17) 具有相容次序。

具有相容次序的矩阵 A 还有一些较重要的性质。为了证明这些性质, 首先假设 $\det(A)$ 中不为零的一般项为

$$m(a) = \pm a_{1\alpha(1)} a_{2\alpha(2)} \cdots a_{n\alpha(n)}$$

对于 $m(a)$ 中的每个元素 $a_{i\alpha(i)}$, 若 $i > \alpha(i)$, 则 $a_{i\alpha(i)}$ 为主对角线下方的元素; 如果 $i < \alpha(i)$, 则 $a_{i\alpha(i)}$ 为主对角线上方的元素。

定理 7.2.7 设 A 是具有相容次序的矩阵, l 和 j 分别为 $m(a)$ 中满足 $i > \alpha(i)$ 和 $i < \alpha(i)$ 的元素的个数, 则 $l = j$, 即具有相容次序的矩阵 A , 其行列式 $\det(A)$ 中非零项中的元素, 取在主对角线上方的个数与下方的个数是相等的。

证明 因 A 具有相容次序, 存在满足定义 7.2.4 的子集 S_1 ,

S_1, \dots, S_l , 即对于 $a_{i\alpha(i)} \neq 0$, 若 $i \in S_k$, 则当 $a(i) > l$ 时, $a(i) \in S_{k+1}$, 而当 $a(i) < l$ 时, $a(i) \in S_{k-1}$. 现在定义一个 n 维向量: $\beta = (\beta_1, \beta_2, \dots, \beta_n)^T$, 其中 β_i 是这样确定的, 当 $i \in S_k$ 时, $\beta_i = k$. 向量 β 具有如下性质: 当 $a_{i\alpha(i)} \neq 0$ 时,

$$\beta_i - \beta_{\alpha(i)} = \begin{cases} 1, & i > \alpha(i); \\ -1, & i < \alpha(i). \end{cases}$$

于是
$$l = \sum_{\substack{i=1 \\ i > \alpha(i)}}^n (\beta_i - \beta_{\alpha(i)}), \quad j = \sum_{\substack{i=1 \\ i < \alpha(i)}}^n (\beta_{\alpha(i)} - \beta_i).$$

因为
$$\begin{aligned} l - j &= \sum_{\substack{i=1 \\ i > \alpha(i)}}^n (\beta_i - \beta_{\alpha(i)}) - \sum_{\substack{i=1 \\ i < \alpha(i)}}^n (\beta_{\alpha(i)} - \beta_i) \\ &= \sum_{i \neq \alpha(i)} (\beta_i - \beta_{\alpha(i)}) + (\beta_i - \beta_i) \\ &= \sum_{i=1}^n (\beta_i - \beta_{\alpha(i)}) = \sum_{i=1}^n \beta_i - \sum_{i=1}^n \beta_{\alpha(i)} = 0, \end{aligned}$$

所以
$$l = j.$$

推论 设 $A \in R^{n \times n}$ 具有相容次序且 $D = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$, 矩阵 E 为 A 的下三角部分的元素组成的下三角矩阵, 矩阵 F 为 A 的上三角部分的元素组成的上三角矩阵, α, k 为任意非零实数, 则

$$\det(\alpha E + \alpha^{-1} F + kD) \quad (7.2.18)$$

的值与 α 无关.

定理 7.2.8 如果 A 具有相容次序且其对角元全不为零, 又矩阵 $B = I - \text{diag}(\alpha_i^{-1}) A = D^{-1} L + D^{-1} U$, 则

- 1) 若 μ 是 B 的 p 重特征值, 则 $-\mu$ 也是 B 的 p 重特征值;
- 2) 若 λ 是 J_ω 的特征值, 则存在 B 的一个特征值 μ , 使

$$\lambda + \omega - 1 = \omega \mu \lambda^{\frac{1}{2}}, \quad (7.2.19)$$

- 3) 如果 μ 是 B 的特征值, λ 满足 (7.2.19), 则 λ 是 J_ω 的一个特征值.

证明 因为 A 具有相容次序, 所以 $B - I$ 也具有相容次序. 根据定理 7.2.7 的推论, 取 $\alpha = \pm 1$, $k = \mu$, 则

$$\det(B - \mu I) = \det(-B - \mu I) = (-1)^n \det(B + \mu I),$$

从而 B 的特征多项式可写为

$$\det(B - \mu I) = (-1)^n \mu^{r-1} (\mu^2 - \mu_1^2) (\mu^2 - \mu_2^2) \cdots (\mu^2 - \mu_{\frac{r}{2}}^2),$$

其中 μ_i 是 B 的非零特征值, 于是结论 1) 成立. 根据 (7.2.12)

$$J_\omega - \lambda I = (I - \omega D^{-1}L)^{-1} [\lambda \omega D^{-1}L + \omega D^{-1}U + (1 - \omega - \lambda)I],$$

所以

$$\det(J_\omega - \lambda I) = \det[\lambda \omega D^{-1}L + \omega D^{-1}U + (1 - \omega - \lambda)I]. \quad (7.2.20)$$

因为 A 具有相容次序且对角元全不为零, 根据定理 7.2.7 的推论可得

$$\begin{aligned} \det(J_\omega - \lambda I) &= (\omega \lambda^{\frac{1}{2}})^n \det(\lambda^{\frac{1}{2}} D^{-1}L + \lambda^{-\frac{1}{2}} D^{-1}U \\ &\quad + \frac{1 - \omega - \lambda}{\omega \lambda^{1/2}} I) \\ &= (\omega \lambda^{\frac{1}{2}})^n \det(D^{-1}L + D^{-1}U \\ &\quad - \frac{\lambda + \omega - 1}{\omega \lambda^{1/2}} I). \end{aligned} \quad (7.2.21)$$

设 λ 是 J_ω 的特征值, 若 $\lambda \neq 0$, 则 $\frac{\lambda + \omega - 1}{\omega \lambda^{1/2}}$ 是 B 的特征值, 记为 μ , 于是

$$\lambda + \omega - 1 = \omega \mu \lambda^{\frac{1}{2}}.$$

如果 $\lambda = 0$, 由 (7.2.20) 得出 $\det(J_\omega) = \det(\omega D^{-1}U + (1 - \omega)I) = 0$, 因而必有 $1 - \omega = 0$, 即 $\omega = 1$. 故对任何 λ 值, (7.2.19) 均成立, 因此结论 2) 成立.

设 μ 是 $B = D^{-1}L + D^{-1}U$ 的特征值, 且 λ 满足 (7.2.19).

若 $\lambda \neq 0$, 则 $\mu = (\lambda + \omega - 1) / \omega \lambda^{\frac{1}{2}}$, 于是由 (7.2.21) 得到

$$\det(J_\omega - \lambda I) = 0$$

即 λ 是 J_ω 的特征值; 如果 $\lambda = 0$, 由 (7.2.19) 得到 $\omega = 1$, 于是由 (7.2.20) 可得 $\det(J_\omega) = 0$, 即 $\lambda = 0$ 是 J_ω 的特征值, 从而结论 3) 成立.

定理 7.2.9 如果 A 是对称正定矩阵, 具有相容次序, 并且 $B = I - D^{-1}A = D^{-1}L + D^{-1}U$, 则

1) B 的特征值均为实数且 $S(B) < 1$

$$2) S(J_\omega) = \begin{cases} \left\{ \frac{\omega S(B) + [\omega^2 S(B)^2 - 4(\omega - 1)]^{\frac{1}{2}}}{2} \right\}^2, & 0 < \omega \leq \bar{\omega} \\ \omega - 1, & \bar{\omega} \leq \omega < 2 \end{cases} \quad (7.2.22)$$

这里 $\bar{\omega}$ 是 SOR 法的最佳松弛因子, 记为

$$\omega_{opt} = \bar{\omega} = \frac{2}{1 + \sqrt{1 - S(B)^2}}, \quad (7.2.23)$$

$$S(J_{\bar{\omega}}) = \bar{\omega} - 1. \quad (7.2.24)$$

证明 由于 A 是对称正定矩阵, 所以 $I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 也为对称矩阵. 又由于

$$B = I - D^{-1}A = D^{-\frac{1}{2}}(I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}})D^{\frac{1}{2}},$$

所以 B 与 $I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 相似, 从而 B 的特征值均为实数.

根据定理 7.2.8, 当 μ 是 B 的特征值时, $-\mu$ 也是 B 的特征值, 而且当 λ 满足

$$\lambda + \omega - 1 = \omega |\mu| \lambda^{\frac{1}{2}}$$

时, λ 是 J_ω 的特征值. 设 $\lambda_1^{\frac{1}{2}}$ 和 $\lambda_2^{\frac{1}{2}}$ 为上式的两个根, 据定理 7.1.1 和定理 7.2.2, 当 $0 < \omega < 2$ 时, $|\lambda_1^{\frac{1}{2}}|$ 和 $|\lambda_2^{\frac{1}{2}}|$ 均小于 1. 由线性代数知识可得

$$|\mu| < 1, \text{ 即 } S(B) < 1.$$

当 μ 是 B 的特征值时, 方程式 (7.2.19) 的两个根

$$\lambda_1^{\frac{1}{2}} = \frac{\omega|\mu| + [\omega^2\mu^2 - 4(\omega - 1)]^{\frac{1}{2}}}{2}, \quad (7.2.24)$$

$$\lambda_2^{\frac{1}{2}} = \frac{\omega|\mu| - [\omega^2\mu^2 - 4(\omega - 1)]^{\frac{1}{2}}}{2} \quad (7.2.25)$$

均为 J_ω 的特征值，而且

当 $\mu = 0$ 时， $\lambda_1 = \lambda_2 = 1 - \omega$ ；

当 $\omega^2\mu^2 - 4(\omega - 1) < 0$ 时， λ_1 与 λ_2 为一对共轭复数且 $|\lambda_1| = |\lambda_2| = |\omega - 1|$ ；

当 $\omega^2\mu^2 - 4(\omega - 1) \geq 0$ 时， λ_1 与 λ_2 是正实数且 $\lambda_1\lambda_2 = (\omega - 1)^2$ 。

因为我们关心的是当 ω 取什么值时 $S(J_\omega)$ 最小，所以只要对 $\omega \in (0, 2)$ 内的每个值，取出 λ_1 与 λ_2 中按模最大的一个来讨论。显然，对于任意实数 μ 和 ω 均有

$$|\lambda_1| \geq |\lambda_2|, \quad |\lambda_1| \geq |\omega - 1|, \quad |\lambda_2| \leq |\omega - 1|. \quad (7.2.26)$$

从(7.2.24)可以看出，对于任意给定的 ω ， $|\lambda_1|$ 是 $|\mu|$ 的单调增函数。因此，取 $\mu = S(B)$ 时， $|\lambda_1|$ 达到最大值 $S(J_\omega)$ ，若记 $S(B) = S_b$ ，则

$$S(J_\omega) = \left| \frac{\omega S_b + [\omega^2 S_b^2 - 4(\omega - 1)]^{1/2}}{2} \right|^2. \quad (7.2.27)$$

(7.2.27)中

$$\omega^2 S_b^2 - 4(\omega - 1) = 0 \quad (7.2.28)$$

的根为

$$\omega_{1,2} = \frac{2(1 \pm \sqrt{1 - S_b^2})}{S_b^2}.$$

因为 $S_b < 1$ ，所以方程(7.2.28)在区间 $(0, 2)$ 内仅有一个根，记为 $\bar{\omega}$ ，则

$$\bar{\omega} = \frac{2(1 - \sqrt{1 - S_b^2})}{S_b^2} = \frac{2}{1 + (1 - S_b^2)^{\frac{1}{2}}}$$

$$= 1 + \left(\frac{S_b}{1 + (1 - S_b^2)^{\frac{1}{2}}} \right)^2. \quad (7.2.29)$$

因此, 当 $0 < \omega < \bar{\omega}$ 时, (7.2.28) 的左边取正值, 于是

$$S(J_\omega) = \left(\frac{\omega S_b + (\omega^2 S_b^2 - 4(\omega - 1))^{\frac{1}{2}}}{2} \right)^2, \quad (7.2.30)$$

当 $\bar{\omega} \leq \omega < 2$ 时, (7.2.28) 的左边取负值, 从而 (7.2.27) 右边 $|\cdot|$ 内取复数, 其模为 $(\omega - 1)^{\frac{1}{2}}$, 则

$$S(J_\omega) = \omega - 1, \quad (7.2.31)$$

因此 (7.2.22) 成立.

(7.2.22) 建立了 $S(J_\omega)$ 与 ω 的函数关系, 现在来求 $S(J_\omega)$ 的最小值. 首先讨论 $0 < \omega < \bar{\omega}$ 的情况, 此时 $S(J_\omega)$ 取 (7.2.22) 的第 1 式, 若令

$$p(\omega) = [\omega^2 S_b^2 - 4(\omega - 1)]^{1/2},$$

则

$$\begin{aligned} p'(\omega) &= \frac{\omega S_b^2 - 2}{p(\omega)}, \\ \frac{d}{d\omega} S(J_\omega) &= \frac{1}{2} (\omega S_b + p(\omega)) (S_b + p'(\omega)) \\ &= \frac{1}{2} (\omega S_b + p(\omega)) \frac{4(S_b^2 - 1)}{p(\omega) [p(\omega) S_b + 2 - \omega S_b^2]} \\ &= \begin{cases} < 0, & 0 < \omega < \bar{\omega} \text{ 时;} \\ -\infty, & \omega \rightarrow \bar{\omega} - 0 \text{ 时.} \end{cases} \end{aligned} \quad (7.2.32)$$

$$\begin{aligned} \frac{d^2}{d\omega^2} S(J_\omega) &= \frac{1}{2} \left[(S_b + p'(\omega))^2 + (\omega S_b + p(\omega)) \frac{S_b^2 - (p'(\omega))^2}{p(\omega)} \right] \\ &= \frac{S_b + p'(\omega)}{2p(\omega)} (2p(\omega) S_b + \omega S_b^2 - \omega S_b p'(\omega)) \\ &= \begin{cases} < 0, & \text{当 } 0 < \omega < \bar{\omega} \text{ 时;} \\ -\infty, & \text{当 } \omega \rightarrow \bar{\omega} - 0 \text{ 时.} \end{cases} \end{aligned}$$

由此可见, 当 $0 < \omega < \bar{\omega}$ 时, $S(J_\omega)$ 是单调下降函数; 当 $\bar{\omega} < \omega < 2$ 时, $S(J_\omega) = \omega - 1$, 是单调上升函数. 因此, 当 $\omega = \bar{\omega}$ 时, $S(J_\omega)$ 达到最小值, 故

$$\omega_{opt} = \bar{\omega} = \frac{2}{1 + \sqrt{1 - (S(B))^2}}.$$

显然, 当 $\omega = \omega_{opt} = \bar{\omega}$ 时, $S(J_\omega) = \bar{\omega} - 1 = \omega_{opt} - 1$.

松弛因子 ω 与 $S(J_\omega)$ 之间的关系, 以及 ω_{opt} 的位置见图 7.1.

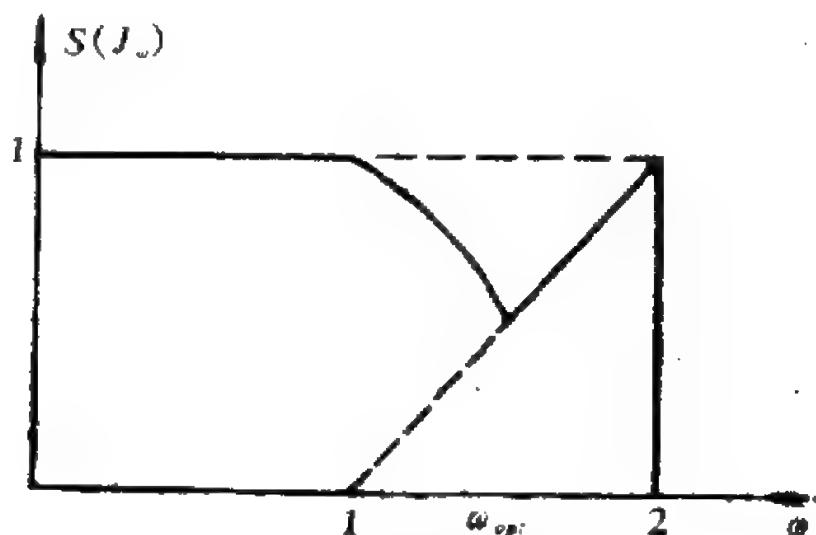


图7.1 超松弛法的谱半径

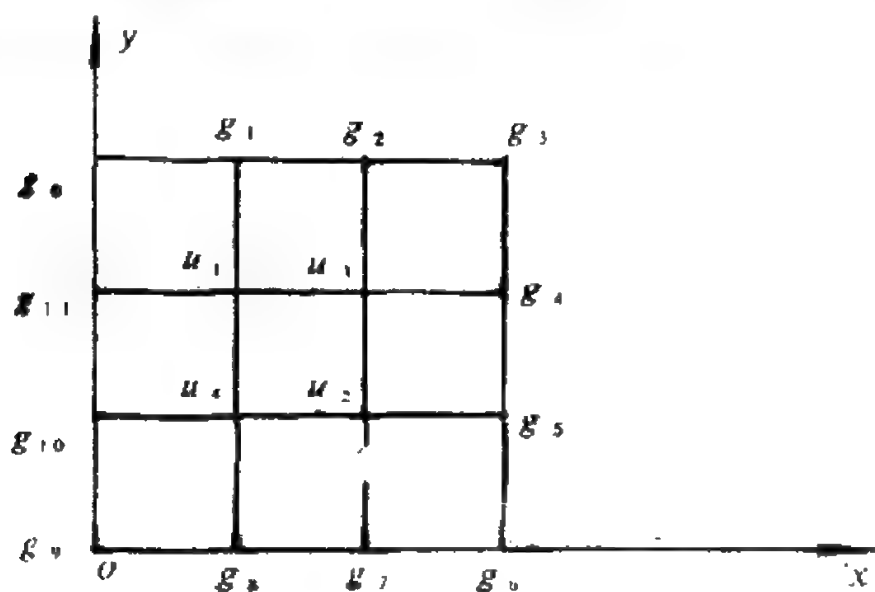
由图 7.1 可以看出, 当 ω 从 ω_{opt} 的左边逐渐增加而趋近于 ω_{opt} 时, $S(J_\omega)$ 的切线斜率趋于 $-\infty$, 即切线趋于垂直于 ω 轴. 但当 ω 从 ω_{opt} 的右边逐渐减小而趋近于 ω_{opt} 时, $S(J_\omega)$ 的切线斜率不变 (恒为 1), 这就是说 $\omega_{opt} - \Delta\omega$ 与 $\omega_{opt} + \Delta\omega$ 中, 后者对应的谱半径将较之前者为小. 所以, 在实际计算 ω_{opt} 的近似值时, 取 ω 略大于 ω_{opt} 比取 ω 略小于 ω_{opt} 更为有利.

例7.2.4.

$$\begin{cases} u_{xx}(x, y) + u_{yy}(x, y) = 0, & 0 < x, y < 1, \\ u|_{\Gamma} = g(x, y) = 0. \end{cases} \quad (7.2.32)$$

其中 Γ 是边长为 1 的正方形边界, 其真解显然为 $u(x, y) \equiv 0$.

在单位正方形上, 用边长 $h = 1/3$ 作出均匀的正方形网格, 并对网格结点进行编号, 如图 7.2 所示.



7.2 网格

记 $u_1 = u\left(\frac{1}{3}, \frac{2}{3}\right)$, $u_2 = u\left(\frac{2}{3}, \frac{1}{3}\right)$, $u_3 = u\left(\frac{2}{3}, \frac{2}{3}\right)$, $u_4 = u\left(\frac{1}{3}, \frac{1}{3}\right)$, 设 u_i 的近似值为 w_i ($i=1, 2, 3, 4$), 求 w_i . 精度要求为 $\|w^{(m)} - w^*\|_\infty \leq 10^{-5}$.

解 将方程 (7.2.32) 离散化后, u_i 的近似值 w_i 可以表示为如下形式:

$$\begin{cases} w_1 = \frac{1}{4}(w_3 + w_4 + g_1 + g_{11}), \\ w_2 = \frac{1}{4}(w_3 + w_4 + g_5 + g_7), \\ w_3 = \frac{1}{4}(w_1 + w_2 + g_2 + g_4), \\ w_4 = \frac{1}{4}(w_1 + w_2 + g_8 + g_{10}). \end{cases} \quad (7.2.33)$$

又因为 $g_i = 0$, 所以 (7.2.33) 可以写成如下形式

$$Aw = 0.$$

其中

$$A = \begin{pmatrix} 1 & 0 & -\frac{1}{4} & -\frac{1}{4} \\ 0 & 1 & -\frac{1}{4} & -\frac{1}{4} \\ -\frac{1}{4} & -\frac{1}{4} & 1 & 0 \\ -\frac{1}{4} & -\frac{1}{4} & 0 & 1 \end{pmatrix}, \quad w = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{pmatrix},$$

$$B = I - D^{-1}A = \frac{1}{4} \begin{pmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{pmatrix}.$$

因为

$$\det(\mu I - B) = \mu^2(\mu^2 - \frac{1}{4}) = 0,$$

所以 B 的特征值为 $\mu_1 = \mu_2 = 0$, $\mu_3 = \frac{1}{2}$, $\mu_4 = -\frac{1}{2}$, 则 $S(B) = \frac{1}{2}$; 又因为 A 的特征值为 $\lambda_1 = \lambda_2 = 1$, $\lambda_3 = 1 - \frac{1}{2} = \frac{1}{2}$, $\lambda_4 = 1 + \frac{1}{2} = 1\frac{1}{2}$, 所以 A 为对称正定矩阵; 据例 7.2.3 知, A 具有相容次序, 故使用 SOR 法的最佳松弛因子 ω_{opt} 可选为

$$\omega_{opt} = \frac{2}{1 + \sqrt{1 - (S(B))^2}} = 1.072.$$

取迭代初值为: $w_1^{(0)} = 0.222222$, $w_2^{(0)} = 0.222222$, $w_3^{(0)} = 4.44444$, $w_4^{(0)} = 0.111111$. 使用 SOR 法计算, 看松弛因子 ω 对迭代次数的影响. 实际迭代次数见下表.

ω	迭 代 次 数	ω	迭 代 次 数
1.020	11	1.072	8
1.040	10	1.202	10
1.052	9	1.402	17

在实际计算中, 由于 $S(B)$ 一般也是难于事先求出的, 于是 ω_{opt} 的值无法事先确定. 因此, 在计算中, 如何确定 ω_{opt} 的近似值是一个比较困难的问题. 这里介绍一种在使用SOR法迭代的过程中顺便求出 ω_{opt} 近似值的方法.

当 $\omega \in (0, \omega_{opt})$ 时, 由(7.2.22)的第1式可以求出 $S(B)$ 用 $S(J_\omega)$ 和 ω 表示的形式

$$S(B) = \frac{S(J_\omega) + (\omega - 1)}{\omega(S(J_\omega))}. \quad (7.2.34)$$

因此, 若能求出 $S(J_\omega)$, $S(B)$ 也就知道了, 同时就能求出 ω_{opt} 的近似值.

首先取一个不太大的正数 ω , 使之满足 $0 < \omega \leq \omega_{opt}$, 然后使用SOR迭代格式

$$x^{(m)} = J_\omega x^{(m-1)} + g$$

进行迭代. 于是有

$$\begin{aligned} x^{(m)} - x^{(m-1)} &= J_\omega (x^{(m-1)} - x^{(m-2)}) = \dots \\ &= J_\omega^{m-1} (x^{(1)} - x^{(0)}). \end{aligned} \quad (7.2.35)$$

当 J_ω 是具有线性初等因子的实矩阵, 并且 $S(J_\omega)$ 是 J_ω 的按模最大的单重特征值时, 由于

$$\frac{\|x^{(m)} - x^{(m-1)}\|_\infty}{\|x^{(m-1)} - x^{(m-2)}\|_\infty} = \lambda^{(m)} \rightarrow S(J_\omega), \text{ 当 } m \rightarrow \infty \text{ 时} \quad (7.2.36)$$

当迭代经过若干步迭代 $\lambda^{(m)}$ 逐渐稳定下来时, 就把 $\lambda^{(m)}$ 作为 $S(J_\omega)$ 的近似值, 再代入(7.2.34)求出 $S(B)$, 然后利用(7.2.23)求出 ω_{opt} 的近似值. 再利用这个 ω_{opt} 代替原来的 ω 继续进行迭代. 直到找到比较满意的 ω_{opt} 的近似值为止. 使用这个方法, 由于难于判断 $S(J_\omega)$ 的近似值是否已较好地接近真值, 以及利用乘幂法求单个最大特征值也要受到有关条件和误差的影响, 要想得到较好的 ω_{opt} 的近似值, 往往需要很多次的迭代.

对于某些特殊问题，矩阵 B 的按模最大的特征值 μ 可以使用某些办法估计出来，这时就可以按 (7.2.23) 算出 ω_{opt} 的近似值，直接用到 SOR 迭代格式中去。但是应该注意，对于具有相容次序的矩阵在求得 ω_{opt} 的近似值后，应该采用比它稍大一点的值，作为实际计算中的松弛因子。

对于那些目前还没有理论公式可以计算最佳松弛因子的矩阵，我们可以用试算的办法来确定它的近似值。最简单的办法是，从同一初始向量出发，用不同的松弛因子 ω 迭代相同的次数，然后比较它们相应的剩余或误差，并选取使剩余或误差之模最小的松弛因子作为最佳松弛因子的近似值。这个方法既简单，又有效。特别，当使用者需要多次求解具有相同系数矩阵的方程组时更是如此。

有时还可以根据求解问题中最佳松弛因子的性质，利用迭代过程中的信息，将它估计出来，而在估算过程中应尽量不要占用过多的计算时间。

§3 预处理和块松弛法

在使用迭代方法之前，特别是当系数阵 A 的条件数很大时，为了加快迭代过程的收敛速度或者为了适应迭代矩阵所具有的性质需要，有时要对原方程组进行事先处理，即预处理。

在使用 SOR 方法时，求新的近似解是按分量逐个地进行的，对于某些方程组，如能适当进行预处理，确定新的近似解则可以成组地进行，有时还能加快收敛的速度。

3.1 预处理法

设线性方程组为

$$Ax = b. \quad (7.3.1)$$

其中 A 为对称正定阵, 为了讨论方便设 A 的主对角线元均为 1, 则 A 能分解成下列形式

$$A = I - L - L^T.$$

于是方程组

$$\begin{cases} B\mathbf{y} = \mathbf{d}, \\ B = (I - \omega L)^{-1} A (I - \omega L)^{-T}, \\ \mathbf{y} = (I - \omega L^T) \mathbf{x}, \quad \mathbf{d} = (I - \omega L)^{-1} \mathbf{b} \end{cases} \quad (7.3.2)$$

与 (7.3.1) 等价。

使用迭代法求解时, 方程组 (7.3.2) 与 (7.3.1) 的迭代矩阵是不同的, 其谱半径一般也不相同, 因此它们的收敛速度也不一样。如果对变形后的方程组 (7.3.2) 使用简单迭代法, 其迭代格式为

$$\mathbf{y}^{(k+1)} = \mathbf{d} + (I - B) \mathbf{y}^{(k)}, \quad (7.3.3)$$

它与原方程组使用简单迭代法计算时, 其步骤及迭代效果都不同。由于 B 是三个矩阵的乘积, 由 (7.3.3) 从 $\mathbf{y}^{(k)}$ 计算 $\mathbf{y}^{(k+1)}$ 时可经下列的代换把计算过程分解为更基本的运算。首先, 令

$$B\mathbf{y}^{(k)} = \mathbf{z}^{(k)}, \quad (7.3.4)$$

则

$$\mathbf{y}^{(k+1)} = \mathbf{d} + \mathbf{y}^{(k)} - \mathbf{z}^{(k)}, \quad (7.3.5)$$

于是

$$\mathbf{z}^{(k)} = (I - \omega L)^{-1} A (I - \omega L)^{-T} \mathbf{y}^{(k)}. \quad (7.3.6)$$

由于

$$(I - \omega L)^{-T} \mathbf{y}^{(k)} = \mathbf{x}^{(k)}, \quad (7.3.7)$$

故

$$(I - \omega L) \mathbf{z}^{(k)} = A \mathbf{x}^{(k)} = \mathbf{F}^{(k)}. \quad (7.3.8)$$

(7.3.4) ~ (7.3.8) 的过程可以用下列算法来描述。

算法 7.3.1. (预处理算法)

- 1) 输入 ω 、 $\mathbf{x}^{(0)}$, 置 $\mathbf{y}^{(0)} = (I - \omega L)^T \mathbf{x}^{(0)}$,
- 2) 由 $(I - \omega L) \mathbf{d} = \mathbf{b}$ 计算 \mathbf{d} ,
- 3) 对 $k = 0, 1, \dots, k_{\max}$ 做
 - 3.1) 置 $\mathbf{F}^{(k)} = A \mathbf{x}^{(k)}$,
 - 3.2) 由 $(I - \omega L) \mathbf{z}^{(k)} = \mathbf{F}^{(k)}$ 计算 $\mathbf{z}^{(k)}$,
 - 3.3) $\mathbf{y}^{(k+1)} = \mathbf{d} + \mathbf{y}^{(k)} - \mathbf{z}^{(k)}$,
 - 3.4) 由 $(I - \omega L)^T \mathbf{x}^{(k+1)} = \mathbf{y}^{(k+1)}$ 计算 $\mathbf{x}^{(k+1)}$,
 - 3.5) 如 $\mathbf{x}^{(k+1)}$ 可接受则转4),
 - 3.6) NEXT k .

4) 结束.

这里 k_{\max} 是允许的最大的迭代次数, ω 是预选的某一适当的因子.

定理7.3.1 若矩阵 A 是主对角线元素为1的对称正定矩阵, 则当 $\omega = 1$ 时, 迭代法(7.3.3)收敛.

证明 要证 (7.3.3) 收敛, 只需证 $S(I - B) < 1$. 设矩阵 B 的特征值为 μ , 其相应的特征向量为 \mathbf{p} , 则

$$B\mathbf{p} = \mu\mathbf{p},$$

即

$$(I - \omega L)^{-1} A (I - \omega L)^{-T} \mathbf{p} = \mu \mathbf{p}.$$

两边用 \mathbf{p}^T 左乘, 得到

$$\mathbf{p}^T (I - \omega L)^{-1} A (I - \omega L)^{-T} \mathbf{p} = \mu \mathbf{p}^T \mathbf{p}, \quad (7.3.9)$$

令

$$\mathbf{v} = (I - \omega L)^{-T} \mathbf{p},$$

由(7.3.9)得到

$$\begin{aligned} \mu &= \frac{\mathbf{v}^T A \mathbf{v}}{\mathbf{p}^T \mathbf{p}} = \frac{\mathbf{v}^T A \mathbf{v}}{\mathbf{v}^T (I - \omega L) (I - \omega L^T) \mathbf{v}} \\ &= \frac{\mathbf{v}^T A \mathbf{v}}{(1 - \omega) \mathbf{v}^T \mathbf{v} + \omega \mathbf{v}^T A \mathbf{v} + \omega^2 \mathbf{v}^T L L^T \mathbf{v}}, \end{aligned} \quad (7.3.10)$$

故当 $\omega = 1$ 时,

$$\mu = \frac{\mathbf{v}^T \mathbf{A} \mathbf{v}}{\mathbf{v}^T \mathbf{A} \mathbf{v} + \mathbf{v}^T \mathbf{L} \mathbf{L}^T \mathbf{v}}. \quad (7.3.11)$$

由于 \mathbf{A} 为对称正定矩阵, 且 $\mathbf{v}^T \mathbf{L} \mathbf{L}^T \mathbf{v} = \mathbf{v}^T \mathbf{L} (\mathbf{v}^T \mathbf{L})^T \geq 0$, 所以 $0 < \mu \leq 1$. 又因矩阵 $\mathbf{I} - \mathbf{B}$ 的特征值为 $1 - \mu$, 所以

$$S(\mathbf{I} - \mathbf{B}) < 1.$$

实际上, ω 的最优值常是大于 1 的, 但这里参数 ω 的取值不象 **SOR** 法那样严格, 通常可用 $\omega = 1$ 作为近似最优值.

对于一般的实对称正定矩阵 \mathbf{A} , 可以将它分解成

$$\mathbf{A} = \mathbf{D} - \bar{\mathbf{L}} - \bar{\mathbf{L}}^T,$$

其中 $\mathbf{D} = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$. 于是方程组

$$\begin{cases} \mathbf{B} \mathbf{y} = \mathbf{d}, \\ \mathbf{B} = (\mathbf{D}^{\frac{1}{2}} - \omega \mathbf{L})^{-1} \mathbf{A} (\mathbf{D}^{\frac{1}{2}} - \omega \mathbf{L})^{-T}, \\ \mathbf{y} = (\mathbf{D}^{\frac{1}{2}} - \omega \mathbf{L}^T) \mathbf{x}, \\ \mathbf{d} = (\mathbf{D}^{\frac{1}{2}} - \omega \mathbf{L})^{-1} \mathbf{b}, \end{cases} \quad (7.3.12)$$

其中 $\mathbf{L} = \bar{\mathbf{L}} \mathbf{D}^{-\frac{1}{2}}$ 与方程组 $\mathbf{A} \mathbf{x} = \mathbf{b}$ 等价. 同样可对 (7.3.12) 使用简单迭代法, 其迭代格式与 (7.3.3) 相同

$$\mathbf{y}^{(k+1)} = \mathbf{d} + (\mathbf{I} - \mathbf{B}) \mathbf{y}^{(k)}.$$

仿算法 7.3.1, 不难写出一般的含有实对称正定矩阵 \mathbf{A} 的方程组 $\mathbf{A} \mathbf{x} = \mathbf{b}$ 的预处理算法.

算法 7.3.1 (一般的实对称正定方程组的预处理算法)

1) 输入 ω , $\mathbf{x}^{(0)}$, 置 $\mathbf{y}^{(0)} = (\mathbf{D}^{\frac{1}{2}} - \omega \mathbf{L}^T) \mathbf{x}^{(0)}$,

2) 由 $(\mathbf{D}^{\frac{1}{2}} - \omega \mathbf{L}) \mathbf{d} = \mathbf{b}$, 计算 \mathbf{d} ;

3) 对 $k = 0, 1, \dots, k_{\max}$ 做

$$3.1) \mathbf{w}^{(k)} = \mathbf{A} \mathbf{x}^{(k)},$$

3.2) 由 $(D^{\frac{1}{2}} - \omega L)x^{(k)} = w^{(k)}$ 计算 $x^{(k)}$,

3.3) $y^{(k+1)} = d + y^{(k)} - x^{(k)}$,

3.4) 由 $(D^{\frac{1}{2}} - \omega L^T)x^{(k+1)} = y^{(k+1)}$ 求 $x^{(k+1)}$,

3.5) 如 $x^{(k+1)}$ 可接受则转4),

3.6) NEXT k ,

4) 结束.

当 A 的条件数 $\text{cond}(A)$ 较大时, 可以通过预处理法来降低条件数. 当 A 是具有“性质 A ”的对称正定矩阵时, 可以证明, 一定存在 ω (如 $\omega = 0$), 使 $\text{cond}(B) \leq \text{cond}(A)$.

如果将 A 进行三角分解

$$A = \bar{L}\bar{L}^T,$$

令 $D^{\frac{1}{2}} - \omega L = \bar{L}$,

则 (7.3.12) 中的矩阵和向量分别换成

$$B = \bar{L}^{-1}A\bar{L}^{-T}, \quad y = \bar{L}^T x, \quad d = \bar{L}^{-1}b.$$

于是 B 的条件数 $\text{cond}(B) = 1$, 此时经一次迭代就得到精确解. 但是求 \bar{L}^{-1} 相当于求 A^{-1} , 因此, 寻找容易求逆的近似三角分解矩阵 \bar{L} , 是一个引人注意的问题. J. A. Meijerink 和 H. A. Van der Vorst 于 1977 年提出的不完全分解法是很有价值的. 将该法用于某些问题的处理, 现已取得了很好的效果.

3.2 块松弛法

分块迭代法在确定新的近似分量时是成组进行的. 先将求解方程组中的方程和未知数进行分组, 使每个方程和每个未知数属于且仅属于一个组. 进行迭代时, 从每个方程组中同时获得的一组近似解的分量, 如此按组进行计算, 直到求得全部新的近似解分量为止.

在建立块迭代格式时, 分组的方法是可以根据需要来选定

的. 下面介绍一种自然数序号分组法. 首先将自然数集 $W = \{1, 2, \dots, n\}$ 分成互不相交的子集: $R_1, R_2, \dots, R_m, R_1 \cup R_2 \cup \dots \cup R_m = W$ 且 $R_i \cap R_j = \emptyset (i \neq j)$. 对于两个分组法, 当且仅当分成的子集的个数和相应子集均相等时, 才称这两个分组法相同. 如果把分组法用 Π_i 表示, 对于 $W = \{1, 2, 3\}$ 的分组就有:

$$\Pi_0: R_1 = \{1\}, \quad R_2 = \{2\}, \quad R_3 = \{3\};$$

$$\Pi_1: R_1 = \{1, 2\}, \quad R_2 = \{3\};$$

$$\Pi_2: R_1 = \{1, 3\}, \quad R_2 = \{2\};$$

$$\Pi_3: R_1 = \{2, 1\}, \quad R_2 = \{3\};$$

$$\Pi_4: R_1 = \{3\}, \quad R_2 = \{1, 2\}.$$

根据分组法的规定, 此时 $\Pi_1 = \Pi_3$, 但 $\Pi_1 \neq \Pi_4$. 通常我们把分组法: $R_k = \{k\}, k = 1, 2, \dots, n$, 叫做 Π_0 .

如果方程组

$$Ax = b$$

经过 Π 分组后, 变成

$$\begin{array}{ccccccc} A_{11} & A_{12} & \cdots & A_{1m} & \bar{x}_1 & = & \bar{b}_1 \\ A_{21} & A_{22} & \cdots & A_{2m} & \bar{x}_2 & = & \bar{b}_2 \\ \cdots & \cdots & \cdots & \cdots & & & \vdots \\ A_{m1} & A_{m2} & \cdots & A_{mm} & \bar{x}_m & = & \bar{b}_m \end{array}, \quad (7.3.12)$$

其中子块 $A_{rs} (r, s = 1, 2, \dots, m)$ 是矩阵 A 中行号属于子集 R_r 中的行, 以及列号属于 R_s 中的列相交叉点的元素按原顺序排列而成的子块. 子列向量 \bar{x}_r 和 \bar{b}_r 分别是列向量 x 和 b 中分量的序号属于 R_r 中的分量按原顺序排列而成的列子向量. 经过这样分块后, 方程组 (7.3.12) 可以写成

$$\sum_{s=1}^m A_{rs} \bar{x}_s = \bar{b}_r \quad (r = 1, 2, \dots, m). \quad (7.3.13)$$

例如, 方程组

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}. \quad (7.3.14)$$

按分组法 Π 分为: $R_1 = \{1, 3\}$, $R_2 = \{2\}$ 时, 方程组 (7.3.14) 可写成 (7.3.13) 的形式为

$$\begin{pmatrix} a_{11} & a_{13} \\ a_{31} & a_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_3 \end{pmatrix} + \begin{pmatrix} a_{12} \\ a_{32} \end{pmatrix} [x_2] = \begin{pmatrix} b_1 \\ b_3 \end{pmatrix}, \quad (7.3.15)$$

$$[a_{21} \ a_{23}] \begin{pmatrix} x_1 \\ x_3 \end{pmatrix} + [a_{22}][x_2] = [b_2].$$

如果方程组 (7.3.13) 中矩阵 $A_{rr} (r = 1, \dots, m)$ 均为非奇异, 则可以仿照简单迭代法的形式写出块简单迭代格式

$$A_{rr} \bar{\mathbf{x}}_r^{(k+1)} = \bar{\mathbf{b}}_r - \sum_{\substack{s=1 \\ s \neq r}}^m A_{rs} \bar{\mathbf{x}}_s^{(k)}, \quad (r = 1, 2, \dots, m) \quad (7.3.16)$$

或

$$\bar{\mathbf{x}}_r^{(k+1)} = \sum_{s=1}^m B_{rs} \bar{\mathbf{x}}_s^{(k)} + \mathbf{c}_r, \quad (r = 1, 2, \dots, m), \quad (7.3.17)$$

其中

$$B_{rs} = \begin{cases} -A_{rs}^{-1} A_{rs}, & \text{当 } r \neq s, \\ 0, & \text{当 } r = s, \end{cases} \quad (7.3.18)$$

$$\mathbf{c}_r = A_{rr}^{-1} \bar{\mathbf{b}}_r. \quad (7.3.19)$$

我们也可以将 (7.3.17) 写成矩阵形式

$$\mathbf{x}^{(k+1)} = B^{(n)} \mathbf{x}^{(k)} + \mathbf{F}^{(n)}. \quad (7.3.20)$$

其中

$$B^{(n)} = (D^{(n)})^{-1} E^{(n)}, \quad \text{而 } \mathbf{F}^{(n)} = (D^{(n)})^{-1} \mathbf{b},$$

$$E^{(n)} = D^{(n)} - A, \quad D^{(n)} = \text{diag}(A_{11}, A_{22}, \dots, A_{mm}).$$

完全仿照逐次松弛法 (7.1.5) 的形式, 同样可以写出块松弛法的计算公式

$$A_{rr} \bar{\mathbf{x}}_r^{(k+1)} = \omega (\bar{\mathbf{b}}_r - \sum_{s=1}^{r-1} A_{rs} \bar{\mathbf{x}}_s^{(k+1)} - \sum_{s=r+1}^m A_{rs} \bar{\mathbf{x}}_s^{(k)})$$

$$= (1 - \omega) A_{rs} \bar{\mathbf{x}}_r^{(k)} \quad (7.3.21)$$

$$(r = 1, 2, \dots, m)$$

或

$$\bar{\mathbf{x}}_r^{(k+1)} = \omega (\mathbf{c}_r - \sum_{s=1}^{r-1} B_{rs} \bar{\mathbf{x}}_s^{(k+1)} - \sum_{s=r+1}^m B_{rs} \bar{\mathbf{x}}_s^{(k)})$$

$$= (1 - \omega) I_r \bar{\mathbf{x}}_r^{(k)}, \quad (7.3.22)$$

$$(r = 1, 2, \dots, m)$$

将(7.3.22)写成矩阵形式, 得

$$\bar{\mathbf{x}}^{(k+1)} = J_\omega^{(n)} \bar{\mathbf{x}}^{(k)} + (I - \omega L^{(n)})^{-1} \omega \mathbf{F}^{(n)}, \quad (7.3.23)$$

其中

$$J_\omega^{(n)} = (I - \omega L^{(n)})^{-1} (\omega U^{(n)} + (1 - \omega) I), \quad (7.3.24)$$

$$L^{(n)} = (D^{(n)})^{-1} A_L^{(n)}, \quad (7.3.25)$$

$$U^{(n)} = (D^{(n)})^{-1} A_U^{(n)}, \quad (7.3.26)$$

而

$$A_L^{(n)} = \begin{pmatrix} 0 & & & & \\ & A_{21} & & & \\ & \vdots & \ddots & & \\ & & & A_{m-1, m-1} & \\ A_{m1} & \cdots & A_{m, m-1} & 0 \end{pmatrix},$$

$$A_U^{(n)} = \begin{pmatrix} 0 & A_{12} & \cdots & A_{1m} \\ & \ddots & \ddots & \vdots \\ & & & A_{m-1, m} \\ & & & & 0 \end{pmatrix}.$$

例如, 设

$$A = \begin{pmatrix} 4 & -1 & 0 & -1 & 0 \\ -1 & 4 & -1 & 0 & -1 \\ 0 & -1 & 4 & 0 & 0 \\ -1 & 0 & 0 & 4 & -1 \\ 0 & -1 & 0 & -1 & 4 \end{pmatrix},$$

$$\Pi_1: R_1 = \{1, 2, 3\}, R_2 = \{4, 5\},$$

则

$$A_{11} = \begin{pmatrix} 4 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 4 \end{pmatrix}, \quad A_{12} = \begin{pmatrix} -1 & 0 \\ 0 & -1 \\ 0 & 0 \end{pmatrix},$$

$$A_{21} = \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \end{bmatrix}, \quad A_{22} = \begin{bmatrix} 4 & -1 \\ -1 & 4 \end{bmatrix},$$

$$D^{(\Pi_1)} = \begin{bmatrix} A_{11} & 0 \\ 0 & A_{22} \end{bmatrix}, \quad A_L^{(\Pi_1)} = \begin{bmatrix} 0 & 0 \\ A_{21} & 0 \end{bmatrix},$$

$$A_U^{(\Pi_1)} = \begin{bmatrix} 0 & A_{12} \\ 0 & 0 \end{bmatrix},$$

$$(D^{(\Pi_1)})^{-1} = \begin{pmatrix} \frac{15}{56} & \frac{4}{56} & \frac{1}{56} & 0 & 0 \\ \frac{4}{56} & \frac{16}{56} & \frac{4}{56} & 0 & 0 \\ \frac{1}{56} & \frac{4}{56} & \frac{15}{56} & 0 & 0 \\ 0 & 0 & 0 & \frac{4}{15} & \frac{1}{15} \\ 0 & 0 & 0 & \frac{1}{15} & \frac{4}{15} \end{pmatrix},$$

$$L^{(\Pi_1)} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ \frac{4}{15} & \frac{1}{15} & 0 & 0 & 0 \\ \frac{1}{15} & \frac{4}{15} & 0 & 0 & 0 \end{pmatrix},$$

$$U^{(n)} = \begin{pmatrix} 0 & 0 & 0 & \frac{16}{56} & \frac{4}{56} \\ 0 & 0 & 0 & \frac{4}{56} & \frac{16}{56} \\ 0 & 0 & 0 & \frac{1}{56} & \frac{4}{56} \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} \cdot$$

利用块松弛法求解，在计算 $\bar{x}_1^{(k+1)}$ 时需要解一个低阶方程组(7.3.21)。通常可以用消去法求解。但是，由于 A_{rr} 不随 k 变化，一个有效的办法是，先把 A_{rr} 进行三角分解，并保存起来，反复利用它来求解(7.3.21)中的 $\bar{x}_r^{(k+1)}$ ($k=1, 2, \dots, m$)。如此所需要完成的计算量就可以减少。

完全仿照定理7.2.2的证明方法，可以证明如下定理。

定理7.3.2 若 A 为实对称正定矩阵，则当 $0 < \omega < 2$ 时，块松弛法(7.3.21)收敛。

对于 n 阶方阵 A ，使用 Π 分组法将 $W = \{1, 2, \dots, n\}$ 分成 m 组后，我们定义 m 阶矩阵 $X = (X_{rs})_{m \times m}$ ，其中

$$X_{rs} = \begin{cases} 0, & \text{当 } A_{rs} = 0; \\ 1, & \text{当 } A_{rs} \neq 0. \end{cases} \quad (7.3.27)$$

定义7.3.1 若矩阵 X 具有“性质A”，则称矩阵 A 具有性质 $A^{(\Pi)}$ 。

定义7.3.2 若矩阵 X 具有相容次序，则称矩阵 A 具有 Π 相容次序。

如果 $\Pi = \Pi_0$ 时 X 有相容次序，那末矩阵 A 就具有通常所称的相容次序。

当矩阵 A 是具有 Π 相容次序的对称正定阵且 $B^{(n)} = I -$

$(D^{(\Pi)})^{-1}A$, 则可以仿照定理 7.2.7 的推论、定理 7.2.8 和定理 7.2.9 的方法证明如下结果:

1) 行列式 $\Delta = \det(\alpha A_L^{(\Pi)} + \alpha^{-1} A_U^{(\Pi)} - k D^{(\Pi)})$ 的值与 α 无关, 其中 $\alpha \neq 0$, k 为任意值.

2) 将 B 换成 $B^{(\Pi)}$, J_ω 换成 $J_\omega^{(\Pi)}$ 后, 有类似于定理 7.2.8 和 $S(B^{(\Pi)}) < 1$ 的结果, 从而可以得到块松弛法的最佳松弛因子为

$$\omega_{opt}^{(\Pi)} = \frac{2}{1 + (1 - (S(B^{(\Pi)}))^2)^{1/2}}, \quad (7.3.28)$$

而且

$$S(J_{\omega_{opt}}^{(\Pi)}(\frac{E}{\omega_{opt}})) = \omega_{opt}^{(\Pi)} - 1. \quad (7.3.29)$$

在实际问题中, 有些问题导出的方程组的系数矩阵 A 往往具有块三对角的形式

$$A = \begin{pmatrix} A_{11} & A_{12} & & \\ & \ddots & \ddots & \\ A_{21} & A_{22} & A_{23} & \\ & \ddots & \ddots & \ddots \\ & & & A_{m-1,m} \\ & & & & A_{mm} \end{pmatrix},$$

这种形式的矩阵显然具有 Π 相容次序. 因此, 只要块三对角阵具有对称正定性, 那末, 逐次松弛法的全部结论都适用于块松弛法.

块松弛法也是一种解线性方程组的有效方法. 特别当 A_{rr} 不太复杂, 而且先对 A_{rr} 进行分解的工作量不太大时, 用该法往往能提高收敛速度. 而且, 对于那些没有“性质 A”但是经适当分块后可以具有性质 $A^{(\Pi)}$ 的矩阵, 使用块松弛法进行求解更为合适. 然而尽管松弛法有时可获得较好的效果, 可是当 A_{rr} 较为复杂或者求解程序效率不高时, 往往并不能使总的计算时间

减少。以上这些问题，在使用块松弛法时也必须考虑，以便根据问题的具体情况决定是否选用块松弛法。

§4 Chebyshev 半迭代法

对于一阶线性定常迭代法 $\mathbf{x}^{(m)} = B\mathbf{x}^{(m-1)} + \mathbf{g}$ ，如果收敛速度很慢，则实际意义并不大。我们希望找到一种方法能加快它的收敛速度。基于一阶线性定常迭代法的 Chebyshev 加速，当迭代矩阵 B 的特征值均为实数时，往往能取得较好的效果。

4.1 一般加速原则

给定一个收敛很慢，甚至不收敛的实数序列： $x_0, x_1, \dots, x_n, \dots$ ，记为 $\{x_i\}$ ，我们可以由它来构造一个新序列： $y_0, y_1, \dots, y_n, \dots$ ，记为 $\{y_i\}$ ，使新序列 $\{y_i\}$ 的收敛比 $\{x_i\}$ 的收敛更快。例如

$$\begin{aligned} y_0 &= x_0, \\ y_1 &= \frac{1}{2}(x_0 + x_1), \\ y_2 &= \frac{1}{3}(x_0 + x_1 + x_2), \\ &\vdots \end{aligned} \quad (7.4.1)$$

可以证明，若 $\{x_i\}$ 收敛，那末 $\{y_i\}$ 一定收敛；甚至当 $\{x_i\}$ 不收敛时， $\{y_i\}$ 也可能收敛。如果 $\{y_i\}$ 收敛，我们称 $\{x_i\}$ 为 Cesàro 可求和。

构造新序列的更一般的原则，可以先考虑下列三角形列阵：

$$\begin{array}{cccc} a_{00} & & & \\ a_{10} & a_{11} & & \\ a_{20} & a_{21} & a_{22} & \\ \vdots & \vdots & \vdots & \vdots \end{array}$$

其中

$$\sum_{k=0}^m a_{mk} = 1, \quad m = 0, 1, 2, \dots \quad (7.4.2)$$

再与 $\{x_i\}$ 生成新序列 $\{y_i\}$, 其中

$$y_i = \sum_{k=0}^i a_{ik} x_k \quad i = 0, 1, 2, \dots \quad (7.4.3)$$

显然, 当对所有的 i 均有 $a_{ik} = 0$ ($k < i$ 时) 而 $a_{ii} = 1$ 时, $\{y_i\}$ 就是 $\{x_i\}$; 当 $a_{ik} = 1/(i+1)$ ($k = 0, 1, \dots, i$) 时, (7.4.3) 就是 (7.4.1).

求解线性方程组 $Ax = b$ 的一阶线性定常迭代法 (7.1.9) 所确定的向量序列 $x^{(0)}, x^{(1)}, \dots$ 收敛于方程组 $Ax = b$ 的解 x^* . 今以以下方式定义新向量序列 $\{u^{(m)}\}$:

$$u^{(m)} = \sum_{k=0}^m a_{mk} x^{(k)}, \quad (7.4.4)$$

其中系数 a_{mk} 满足条件 (7.4.2). 由此所定义的过程称为关于一阶线性定常迭代法 (7.1.9) 的一种半迭代法.

在条件 (7.4.2) 的限制下, 我们总可以找到适当的系数, 使 $\{u^{(m)}\}$ 收敛于方程组 $Ax = b$ 的解 x^* . 为了分析半迭代法的收敛性, 令

$$\begin{cases} \varepsilon^{(m)} = x^{(m)} - x^*, \\ \eta^{(m)} = u^{(m)} - x^*. \end{cases} \quad (7.4.5)$$

从 (7.4.4) 和 (7.4.2) 可以推得

$$\begin{aligned} \eta^{(m)} &= u^{(m)} - x^* = \sum_{k=0}^m a_{mk} x^{(k)} - x^* = \sum_{k=0}^m a_{mk} \varepsilon^{(k)}, \\ &= \sum_{k=0}^m a_{mk} (x^{(k)} - x^*) = \sum_{k=0}^m a_{mk} \varepsilon^{(k)}. \end{aligned}$$

由 (7.1.12) 得到

$$\varepsilon^{(k)} = B^k \varepsilon^{(0)} = B^k \eta^{(0)},$$

因此

$$\eta^{(m)} = \sum_{k=0}^m \alpha_{mk} B^k \varepsilon^{(0)} = p_m(B) \eta^{(0)}, \quad (7.4.6)$$

其中

$$p_m(B) = \sum_{k=0}^m \alpha_{mk} B^k \text{ 且 } p_m(1) = 1. \quad (7.4.7)$$

为了使 $\{u^{(m)}\}$ 收敛于 x^* , 即 $\eta^{(m)} \rightarrow 0$, 就需要 $\lim_{m \rightarrow \infty} p_m(B) = 0$. 满足这些要求的多项式 $p_m(x)$ 是存在的, 例如 $p_m(x) = x^m$, 此时 $p_m(B) = B^m$, 由于 $S(B) < 1$, 自然就有 $\lim_{m \rightarrow \infty} p_m(B) = 0$, 且 $p_m(1) = 1^m = 1$, 即 $u^{(m)} \rightarrow x^*$. 那末, 应该选取怎样的多项式 (7.4.7) 才能使向量序列 $\{u^{(m)}\}$ 的收敛最快? 为此, 必须首先弄清收敛快慢的度量.

类似于 §1, 我们首先假定 B 具有线性无关的特征向量系 x_1, x_2, \dots, x_n , 而 $\lambda_1, \lambda_2, \dots, \lambda_n$ 是相应的特征值, 于是 $\eta^{(0)}$ (即 $\varepsilon^{(0)}$) 可以表示为

$$\begin{aligned} \eta^{(0)} &= a_1 x_1 + a_2 x_2 + \dots + a_n x_n, \\ \eta^{(m)} &= p_m(B) \eta^{(0)} = a_1 p_m(B) x_1 + \dots + a_n p_m(B) x_n \\ &= p_m(\lambda_1) a_1 x_1 + p_m(\lambda_2) a_2 x_2 + \dots + p_m(\lambda_n) a_n x_n. \end{aligned} \quad (7.4.8)$$

由此可以看出, $\eta^{(m)}$ 收敛于零的速度是由 $\max_{1 \leq i \leq n} |p_m(\lambda_i)|$ 决定的. 为使 $\{u^{(m)}\}$ 收敛最快, 只需选取所有 m 阶多项式 $p_m(x)$ 中使下式达到极小值者

$$\begin{cases} S(p_m(B)) = \max_{1 \leq i \leq n} |p_m(\lambda_i)|, \\ p_m(1) = 1, \end{cases}$$

因而, 可以类似于定义 7.1.2 一样定义渐近收敛速度.

定义 7.4.1 称 $R_\infty(p_m(B)) = \lim_{m \rightarrow \infty} (-\frac{1}{m} \ln(S(p_m(B))))$ 为半迭代法 (7.4.4) 的渐近收敛速度.

当 $p_m(x) = x^m$ 时, $R_m(p_m(B)) = R(B)$, 因此, 定义 7.4.1 是定义 7.1.2 的自然推广.

4.2 Chebyshev 半迭代法

假定 B 的全部特征值 λ_i 都是实数, 并且满足如下条件

$$\alpha \leq \lambda_i \leq \beta < 1 \quad (\alpha < \beta), \quad (7.4.9)$$

其中 α 可以小于 -1 . 在这种条件下, 用 Chebyshev 多项式构成半迭代法能取得加快收敛速度的效果. 通常将这样的半迭代法称为 Chebyshev 半迭代法. 由于 λ_i 事先不知道, 所以 $\max_{1 \leq i \leq n} |p_m(\lambda_i)|$ 的极小值也难于求出. 通常用求

$$\begin{cases} \max_{\alpha < \lambda < \beta} |p_m(\lambda)|, \\ p_m(1) = 1 \end{cases} \quad (7.4.10)$$

的极小化问题来代替. 当然有

$$\max_{1 \leq i \leq n} |p_m(\lambda_i)|^{\frac{1}{m}} \leq \max_{\alpha < \lambda < \beta} |p_m(\lambda)|^{\frac{1}{m}}. \quad (7.4.11)$$

今引入新的变量

$$\gamma = \gamma(\lambda) = \frac{2\lambda - (\alpha + \beta)}{\beta - \alpha}. \quad (7.4.12)$$

显然有 $\gamma(\alpha) = -1$, $\gamma(\beta) = 1$, 而且当 λ 满足 (7.4.9) 时, $-1 \leq \gamma \leq 1$. 令

$$z = \gamma(1) = \frac{2 - (\alpha + \beta)}{\beta - \alpha}. \quad (7.4.13)$$

又因为

$$2 - (\alpha + \beta) > 2 - (2\beta) = 2(1 - \beta) > 0.$$

$$\beta - \alpha > 0.$$

$$2 - (\alpha + \beta) - (\beta - \alpha) = 2 - 2\beta > 0.$$

所以

$$z > 1.$$

由(7.4.12)解出 $\lambda = \frac{1}{2}[(\beta - \alpha)\gamma + \beta + \alpha]$ 代入 $p_m(\lambda)$ 得到

$$p_m(\lambda) = p_m\left(\frac{(\beta - \alpha)\gamma + \beta + \alpha}{2}\right) = Q_m(\gamma). \quad (7.4.14)$$

因此

$$\max_{\alpha \leq \lambda \leq \beta} |p_m(\lambda)| = \max_{-1 \leq \gamma \leq 1} |Q_m(\gamma)|$$

于是极小化问题(7.4.10)变成

$$\begin{cases} \max_{-1 \leq \gamma \leq 1} |Q_m(\gamma)| = \min \\ Q_m(z) = 1 \end{cases} \quad (7.4.15)$$

这个问题的解答是

$$Q_m(\gamma) = T_m(\gamma) / T_m(z). \quad (7.4.16)$$

其中 $T_m(\gamma)$ 就是 m 阶 Chebyshev 多项式, 其参数表达式可以写成

$$\begin{cases} x = \cos \theta, \\ T_m(x) = \cos m\theta. \end{cases} \quad 0 \leq \theta \leq \pi.$$

这里

$$\begin{aligned} T_m(x) &= \cos(m \arccos x) \\ &= \frac{(x + \sqrt{x^2 - 1})^m + (x - \sqrt{x^2 - 1})^m}{2}, \end{aligned} \quad (7.4.17)$$

它对 $x > 1$ 时也是有意义的. 由(7.4.15)的解可得(7.4.10)的解为

$$\begin{aligned} p_m(\lambda) &= Q_m(\gamma) = Q_m\left(\frac{2\lambda - (\beta + \alpha)}{\beta - \alpha}\right) \\ &= T_m\left(\frac{2\lambda - (\beta + \alpha)}{\beta - \alpha}\right) / T_m\left(\frac{2 - (\beta + \alpha)}{\beta - \alpha}\right), \end{aligned} \quad (7.4.18)$$

而且

$$\begin{aligned} \min \left(\max_{\alpha \leq \lambda \leq \beta} |p_m(\lambda)| \right) &= \max_{\alpha \leq \lambda \leq \beta} \left| T_m\left(\frac{2\lambda - (\beta + \alpha)}{\beta - \alpha}\right) / T_m(z) \right| \\ &= \frac{1}{|T_m(z)|} \end{aligned} \quad (7.4.19)$$

有了 $p_m(\lambda)$ 的表达式(7.4.18)后,利用 Chebyshev 多项式 $T_m(x)$ 的递推性质

$$\begin{cases} T_{m+1}(x) = 2xT_m(x) - T_{m-1}(x) & (m = 1, 2, \dots), \\ T_0(x) = 1, \quad T_1(x) = x, \end{cases} \quad (7.4.20)$$

再通过计算 $p_m(\lambda) = \alpha_{m0} + \alpha_{m1}\lambda' + \dots + \alpha_{mm}\lambda^m$ 的系数 α_{mi} 就可以构造 $u^{(m)}$. 因为

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_2(x) = 2x^2 - 1,$$

可以求得

$$p_0(\lambda) = \alpha_{00} = 1, \quad u^{(0)} = x^{(0)};$$

$$\begin{aligned} p_1(\lambda) &= \alpha_{10} + \alpha_{11}\lambda = T_1\left(\frac{2\lambda - (\beta + \alpha)}{\beta - \alpha}\right) / T_1(z) \\ &= \frac{2\lambda - (\beta + \alpha)}{\beta - \alpha} \cdot \frac{\beta - \alpha}{2 - (\beta + \alpha)} = \frac{2\lambda - (\beta + \alpha)}{2 - (\beta + \alpha)}. \end{aligned}$$

经比较系数, 得

$$\alpha_{10} = -\frac{\beta + \alpha}{2 - (\beta + \alpha)}, \quad \alpha_{11} = \frac{2}{2 - (\beta + \alpha)}.$$

于是

$$\begin{aligned} u^{(1)} &= \frac{-(\beta + \alpha)}{2 - (\beta + \alpha)} x^{(0)} + \frac{2}{2 - (\beta + \alpha)} x^{(1)} \\ &= \frac{2}{2 - (\beta + \alpha)} (Bx^{(0)} + \mathbf{g}) - \frac{(\beta + \alpha)}{2 - (\beta + \alpha)} x^{(0)}. \end{aligned} \quad (7.4.21)$$

如此继续进行, 就可以求出 $u^{(m)}$ 来。但是, 这种算法比较麻烦, 下面我们导出利用 $u^{(m-1)}$, $u^{(m)}$ 来计算 $u^{(m+1)}$ 的递推关系式。从(7.4.18)可得

$$p_m(B) = T_m\left(\frac{2B - (\beta + \alpha)I}{\beta - \alpha}\right) / T_m(z),$$

由(7.4.6)和(7.4.20)可以推出

$$\begin{aligned}
\eta^{(m+1)} &= \left[T_{m+1} \left(\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right) / T_{m+1}(z) \right] \varepsilon^{(0)} \\
&= \left\{ \left[2 \left(\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right) T_m \left(\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right) \right. \right. \\
&\quad \left. \left. - T_{m-1} \left(\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right) \right] \cdot \left[T_{m+1}(z) \right]^{-1} \right\} \varepsilon^{(0)} \\
&= 2 \left[\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right] \frac{T_m(z)}{T_{m+1}(z)} \eta^{(m)} \\
&\quad - \frac{T_{m-1}(z)}{T_{m+1}(z)} \eta^{(m-1)}, \tag{7.4.22}
\end{aligned}$$

将 $\eta^{(m)} = u^{(m)} - x^*$ 代入上式后得到

$$\begin{aligned}
u^{(m+1)} &= 2 \left[\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right] \frac{T_m(z)}{T_{m+1}(z)} u^{(m)} \\
&\quad - \frac{T_{m-1}(z)}{T_{m+1}(z)} u^{(m-1)} + \left[I - 2 \left(\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right) \right. \\
&\quad \left. \cdot \frac{T_m(z)}{T_{m+1}(z)} + \frac{T_{m-1}(z)}{T_{m+1}(z)} I \right] x^*. \tag{7.4.23}
\end{aligned}$$

但是

$$\begin{aligned}
&\left[I - 2 \left(\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right) \frac{T_m(z)}{T_{m+1}(z)} + \frac{T_{m-1}(z)}{T_{m+1}(z)} I \right] x^* \\
&= 2 \left[zI - \frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right] \frac{T_m(z)}{T_{m+1}(z)} x^* \\
&= \frac{4(I - B)}{\beta - \alpha} \frac{T_m(z)}{T_{m+1}(z)} x^* = \frac{4}{\beta - \alpha} \frac{T_m(z)}{T_{m+1}(z)} \mathbf{E}, \tag{7.4.24}
\end{aligned}$$

因此, 当 $m \geq 1$ 时,

$$u^{(m+1)} = 2 \left[\frac{2B - (\beta + \alpha)I}{\beta - \alpha} \right] \frac{T_m(z)}{T_{m+1}(z)} u^{(m)} - \frac{T_{m-1}(z)}{T_{m+1}(z)} u^{(m-1)} \\ + \frac{4}{\beta - \alpha} \frac{T_m(z)}{T_{m+1}(z)} g, \quad (7.4.25)$$

由(7.4.21)和(7.4.25)可以求出 $\{u^{(m)}\}$ 。如果令 $\rho_1 = 1$,

$$\rho_m = 2zT_{m-1}(z)/T_m(z), \quad m = 2, 3, \dots \quad (7.4.26)$$

(7.4.25)还可以简化成

$$u^{(m+1)} = \frac{\rho_{m+1}}{2 - (\beta + \alpha)} \left\{ [2B - (\beta + \alpha)I] u^{(m)} + 2g \right\} \\ + (1 - \rho_{m+1}) u^{(m-1)}.$$

将 $T_{m+1}(x) = 2xT_m(x) - T_{m-1}(x)$ 的两边乘上 $1/2xT_m(x)$ 得到

$$\frac{T_{m+1}(x)}{2xT_m(x)} = 1 - \frac{T_{m-1}(x)}{2xT_m(x)} = 1 - \frac{1}{4x^2} \frac{2xT_{m-1}(x)}{T_m(x)}.$$

由(7.4.26)可得

$$\rho_{m+1}^{-1} = 1 - \frac{1}{4z^2} \rho_m,$$

而

$$\rho_2 = \frac{2zT_1(z)}{T_2(z)} = \frac{2z^2}{2z^2 - 1},$$

从而得到基于一阶线性定常迭代法(7.1.9)的 Chebyshev 迭代法的简化计算公式为

$$\begin{aligned} & \text{初始值 } x^{(0)} = u^{(0)}, \\ & u^{(1)} = \frac{2}{2 - (\alpha + \beta)} (Bx^{(0)} + g) - \frac{\alpha + \beta}{2 - (\alpha + \beta)} x^{(0)}, \\ & u^{(m+1)} = \frac{\rho_{m+1}}{2 - (\alpha + \beta)} \{ [2B - (\alpha + \beta)I] u^{(m)} + 2g \} \\ & \quad + (1 - \rho_{m+1}) u^{(m-1)}. \end{aligned} \quad (7.4.27)$$

$$\rho_1 = 1, \quad \rho_2 = \frac{2z^2}{2z^2 + 1},$$

其中,
$$\rho_{m+1} = \left[1 - \frac{1}{4z} \rho_m \right]^{-1}, \quad m = 2, 3, \dots$$

$$z = \frac{2}{\beta - \alpha} \frac{(\alpha + \beta)}{2}.$$

从这些计算公式中可以看出, **Chebyshev** 半迭代法是一种二阶线性非定常迭代法. 使用该迭代法时要注意参数 α 和 β 的选取, 在一般情况下可运用特征值估计的有关方法来定 α 和 β .

下面我们来讨论 **Chebyshev** 半迭代法的收敛性.

由(7.4.17)可得

$$\begin{aligned} T_m(z) &= \frac{(z + \sqrt{z^2 - 1})^m - (z - \sqrt{z^2 - 1})^m}{2} \\ &= \frac{(1 + \sqrt{1 - 1/z^2})^m + (1 - \sqrt{1 - 1/z^2})^m}{(2/z^m)}. \end{aligned}$$

令 $\delta = 1/z$, 则

$$\begin{aligned} T_m(z) &= \frac{(1 + \sqrt{1 - \delta^2})^m + (1 - \sqrt{1 - \delta^2})^m}{2\delta^m} \\ &= \frac{1}{2} \left[\left(\frac{1 + \sqrt{1 - \delta^2}}{\delta} \right)^m + \left(\frac{1 - \sqrt{1 - \delta^2}}{\delta} \right)^m \right]. \end{aligned}$$

令 $\tau = \delta / (1 + \sqrt{1 - \delta^2})$, 则

$$T_m(z) = \frac{1}{2} \left[\frac{1}{\tau^m} + \tau^m \right] = \frac{1 + \tau^{2m}}{2\tau^m},$$

从而推得

$$\text{当 } z > 1 \text{ 时, } [T_m(z)]^{-1} = \frac{2\tau^m}{1 + \tau^{2m}},$$

其中

$$\tau = \frac{\delta}{1 + \sqrt{1 - \delta^2}}, \quad \delta = \frac{1}{z}.$$

显然, 当 $m \rightarrow \infty$ 时, $[T_m(z)]^{-1} \rightarrow 0$. 于是根据(7.4.8)和(7.4.19)可以推出

当 $m \rightarrow \infty$ 时, $\eta^{(m)} = u^{(m)} - x^* \rightarrow 0$,

即 $u^{(m)} \rightarrow x^*$. 因此, 不管一阶线性定常迭代法 (7.1.9) 是否收敛, 对于满足上述条件的 Chebyshev 半迭代法是一定收敛的.

如果原迭代法收敛, 它的渐近收敛速度为 $R(B)$, Chebyshev 半迭代法的渐近收敛速度为 $R(p_m(B))$, 那末, 两者收敛速度的快慢如何呢? 下面将作出比较.

据定义 7.4.1, $R(p_m(B)) = \lim_{m \rightarrow \infty} (-\frac{1}{m} \ln S(p_m(B)))$, 由 (7.4.11) 可以推出

$$S(p_m(B)) \leq (T_m(z))^{-1}$$

从而

$$\begin{aligned} R(p_m(B)) &\geq \lim_{m \rightarrow \infty} \left(-\frac{1}{m} \ln \frac{1}{T_m(z)} \right) \\ &= \lim_{m \rightarrow \infty} \left(-\frac{1}{m} \ln \frac{2\tau^m}{1 + \tau^{2m}} \right) = -\ln \tau. \end{aligned} \quad (7.4.28)$$

另一方面, 我们可以作出线性函数 $\bar{p}_1(x) = a_1 x + a_0$, 其中 $a_0 + a_1 = 1$, 使

$$\max_{\alpha \leq \lambda \leq \beta} |\bar{p}_1(\lambda)| \leq \max_{1 \leq i \leq n} |\lambda_i|.$$

这里 λ_i 和 $\bar{p}_1(\lambda_i)$ 分别是 B 和 $\bar{p}_1(B)$ 的特征值. 由此可得

$$R(B) \leq R(\bar{p}_1(B)) \leq -\ln \frac{1}{T_1(z)} = -\ln \sigma.$$

因此,

$$\frac{R(p_m(x))}{[R(B)]^{1/2}} \geq \frac{-\ln \tau}{[-\ln \sigma]^{1/2}},$$

其中 $\tau = \frac{\sigma}{1 + \sqrt{1 - \sigma^2}}$, $0 < \sigma < 1$.

现在我们来证明

$$\lim_{\sigma \rightarrow 1} \frac{-\ln \tau}{[-\ln \sigma]^{1/2}} = \sqrt{2}.$$

令 $\sigma = e^{-a}$, 则 $a = -\ln\sigma$, 由于

$$\begin{aligned} -\ln \frac{\sigma}{1 + \sqrt{1 - \sigma^2}} &= -\frac{1}{2} \ln \frac{1 - \sqrt{1 - \sigma^2}}{1 + \sqrt{1 - \sigma^2}} \\ &\geq (1 - \sigma^2)^{\frac{1}{2}} = (1 - e^{-2a})^{\frac{1}{2}} \\ &= e^{-a} (e^{2a} - 1)^{\frac{1}{2}} \geq \sigma (2a)^{\frac{1}{2}} \\ &= \sigma (-2\ln\sigma)^{\frac{1}{2}}, \end{aligned}$$

并且

$$\begin{aligned} -\ln \frac{\sigma}{1 + \sqrt{1 - \sigma^2}} &= a + \ln(1 + \sqrt{1 - \sigma^2}) \leq a + (1 - \sigma^2)^{\frac{1}{2}} \\ &= a + (1 - e^{-2a})^{\frac{1}{2}} \leq a + (1 - (1 - 2a))^{\frac{1}{2}} \\ &= -\ln\sigma + (-2\ln\sigma)^{\frac{1}{2}}, \end{aligned}$$

所以,

$$\frac{\sigma (-2\ln\sigma)^{\frac{1}{2}}}{(-2\ln\sigma)^{\frac{1}{2}}} \leq \frac{-\ln \frac{\sigma}{1 + \sqrt{1 - \sigma^2}}}{(-2\ln\sigma)^{\frac{1}{2}}} \leq \frac{-\ln\sigma + (-2\ln\sigma)^{\frac{1}{2}}}{(-2\ln\sigma)^{\frac{1}{2}}}.$$

由于 $\sigma \rightarrow 1^-$ 时,

$$\frac{\sigma (-2\ln\sigma)^{\frac{1}{2}}}{(-2\ln\sigma)^{\frac{1}{2}}} \rightarrow 1, \quad \frac{-\ln\sigma + (-2\ln\sigma)^{\frac{1}{2}}}{(-2\ln\sigma)^{\frac{1}{2}}} \rightarrow 1,$$

故

$$\lim_{\sigma \rightarrow 1^-} \frac{-\ln \frac{\sigma}{1 + \sqrt{1 - \sigma^2}}}{(-2\ln\sigma)^{\frac{1}{2}}} = 1,$$

即

$$\lim_{\sigma \rightarrow 1^-} \frac{-\ln\tau}{(-\ln\sigma)^{\frac{1}{2}}} = \sqrt{2}.$$

这就是说, 当 σ 接近于 1 时, Chebyshev 半迭法的渐近收敛速度比原来的一阶线性定常迭代法 (7.1.9) 的渐近收敛速度

快得多。因此，当 B 的特征值全部是实数时，采用 Chebyshev 半迭代法来加速迭代过程是有效的。

4.3 简单迭代法的 Chebyshev 加速

前面我们讨论了一般的一阶线性定常迭代法的 Chebyshev 加速问题。如果满足条件的某种迭代法已经给出，问题就转化为如何选择其中的参数了。我们以简单迭代法为例来说明这个问题。

设线性方程组为 $Ax = b$ ，其系数矩阵 A 是对称正定的，其相应的简单迭代法为：

$$x^{(m+1)} = Gx^{(m)} + c, \quad (7.4.29)$$

其中 $G = I - D^{-1}A$, $A = \text{diag}(a_{ii})$, $c = D^{-1}b$,
因为

$$\bar{G} = D^{\frac{1}{2}}GD^{-\frac{1}{2}} = I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$$

且 $D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 是对称正定阵，所以 G 的特征值 μ 均为实数且小于 1，于是存在实数 α, β 使

$$\alpha \leq \mu \leq \beta < 1.$$

据 (7.4.27)，简单迭代法的 Chebyshev 半迭法可以表为

$$u^{(m+1)} = \frac{\rho_{m+1}}{2 - (\alpha + \beta)} \left\{ [2G - (\beta + \alpha)I]u^{(m)} + 2c \right\} + (1 - \rho_{m+1})u^{(m-1)}. \quad (7.4.30)$$

在 (7.4.30) 中， G 的特征值进一步满足某些特殊条件时，计算格式变得更为简单些。

1) 当 G 的所有特征值均为小于 1 的相等实数，即 $\alpha = \beta$ 时，显然，可取 $\rho_m = 1$, ($m = 1, 2, \dots$) 从而 (7.4.30) 变成

$$u^{(m+1)} = \frac{1}{1 - \alpha} [(G - \alpha I)u^{(m)} + c]. \quad (7.4.31)$$

又因为 $S(G - \alpha I) = 0$ ，所以迭代格式 (7.4.31) 的收敛速度是非常快的。

2) 如果矩阵 A 是具有相容次序的对称正定阵, 此时 $S(G) < 1$, 故可设 $-\alpha - \beta = S(G)$, 于是

$$z = \frac{2 - (\alpha + \beta)}{\beta - \alpha} = \frac{2}{2S(G)} = \frac{1}{S(G)}. \quad (7.4.32)$$

将(7.4.32)代入(7.4.27)得到

$$\begin{aligned} \rho_1 &= 1, & \rho_2 &= 2/(2 - (S(G))^2), \\ \rho_{m+1} &= (1 - \frac{1}{4}(S(G))^2\rho_m)^{-1} \quad (m = 2, 3, \dots), \end{aligned} \quad (7.4.33)$$

而(7.4.30)变成

$$u^{(m+1)} = \rho_{m+1}(Gu^{(m)} + c) + (1 - \rho_{m+1})u^{(m-1)}. \quad (7.4.34)$$

因此, 在这种条件下, **Chebyshev** 半迭代法的计算公式可以由(7.4.34)和(7.4.33)表示.

根据上面的讨论和(7.4.28), 半迭代法的渐近收敛速度应满足

$$R(p_m(G)) \geq -\ln \tau, \quad \tau = \frac{S(G)}{1 + \sqrt{1 - (S(G))^2}},$$

即

$$\begin{aligned} R(p_m(G)) &\geq -\ln \frac{S(G)}{1 + \sqrt{1 - (S(G))^2}} \\ &= -\frac{1}{2} \ln \frac{1 - \sqrt{1 - (S(G))^2}}{1 + \sqrt{1 - (S(G))^2}}. \end{aligned} \quad (7.4.35)$$

如果改用逐次松弛法, 并取最佳松弛因子为 ω_{opt} , 其渐近收敛速度为

$$\begin{aligned} R(J_{\omega_{opt}}) &= -\ln(\omega_{opt} - 1) = -\ln \left(\frac{2}{1 + \sqrt{1 - (S(G))^2}} - 1 \right) \\ &= -\ln \frac{1 - \sqrt{1 - (S(G))^2}}{1 + \sqrt{1 - (S(G))^2}}. \end{aligned} \quad (7.4.36)$$

当矩阵 A 是具有相容次序的对称正定矩阵时, 比较(7.4.35)和(7.4.36), 可以看出, **SOR** 法的渐近收敛速度为基于简单迭代法的 **Chebyshev** 半迭代法的两倍, 因此, 当线性方程组的系数矩阵 A 是具有相容次序的对称正定阵时, 一般说来, 采取逐次松弛法是较为有利的.

§5 非线性迭代法

非线性迭代法的特点是新近似解是已有近似解的非线性函数. 本节只介绍解线性方程组的最速下降法和共轭斜量法等两种非线性迭代法. 它适用于系数矩阵为对称正定阵的情况, 由于该法不需要选取参数, 所以使用比较方便. 共轭斜量法, 实质上是一种直接法, 如果没有舍入误差, 对于 n 阶方程组最多迭代 n 步便可求得精确解. 但是共轭斜量法又具有迭代的计算公式, 所以就其形式来说又是一种迭代法. 它还可以直接应用于解非线性方程组和求函数的极值问题, 因此, 近年来受到广泛的重视. 它与其它迭代法相比较, 不足的是, 存储量较大和计算过程中对舍入误差比较敏感.

5.1 线性方程组的等价问题

设有线性方程组

$$\sum_{j=1}^n a_{ij}x_j = b_i \quad (i = 1, 2, \dots, n) \quad (7.5.1)$$

或 $Ax = b$, 其中系数矩阵 A 是对称正定的.

考虑下列二次函数的极小值问题

$$\begin{aligned} F(x) &= \frac{1}{2} \sum_{i,j=1}^n a_{ij}x_ix_j - \sum_{i=1}^n b_ix_i \\ &= \frac{1}{2} (Ax, x) - (b, x). \end{aligned} \quad (7.5.2)$$

这里 A 仍为对称正定矩阵, $\mathbf{b} = (b_1, b_2, \dots, b_n)^T$, 如果我们把 n 维向量看成是 n 维空间 R^n 中的一点, 那末, 我们可以证明下列定理.

定理 7.5.1 向量 \mathbf{x}^* 为线性方程组 (7.5.1) 的解的充要条件是, \mathbf{x}^* 是极小值问题 (7.5.2) 的极小点. 即

$$F(\mathbf{x}^*) = \min_{\mathbf{x} \in R^n} F(\mathbf{x}).$$

证明 因为 $a_{ij} = a_{ji}$, 所以

$$\frac{\partial F}{\partial x_i} = \sum_{j=1}^n a_{ij}x_j - b_i = -r_i,$$

$$\text{即} \quad \text{grad} F(\mathbf{x}) = A\mathbf{x} - \mathbf{b} = -\mathbf{r}. \quad (7.5.3)$$

也就是说, 将任意向量 \mathbf{x} 代入方程 (7.5.1) 中所得到的剩余向量为 $\mathbf{r} = \mathbf{b} - A\mathbf{x}$, 则 $-\mathbf{r}$ 就是函数 $F(\mathbf{x})$ 在 \mathbf{x} 处的梯度向量. 若

\mathbf{x}^* 使函数 $F(\mathbf{x})$ 取极小值, 即 $F(\mathbf{x}^*) = \min_{\mathbf{x} \in R^n} F(\mathbf{x})$, 则

$$\text{grad} F(\mathbf{x})|_{\mathbf{x}^*} = A\mathbf{x}^* - \mathbf{b} = \mathbf{0},$$

即 \mathbf{x}^* 是线性方程组 (7.5.1) 的解.

反之, 若 \mathbf{x}^* 是线性方程组 (7.5.1) 的解, 而 A 又是对称正定阵, 故有

$$\begin{aligned} F(\mathbf{x}) - F(\mathbf{x}^*) &= -\frac{1}{2} (A\mathbf{x}, \mathbf{x}) - (\mathbf{b}, \mathbf{x}) - \left[-\frac{1}{2} (A\mathbf{x}^*, \mathbf{x}^*) \right. \\ &\quad \left. + (\mathbf{b}, \mathbf{x}^*) \right] \\ &= -\frac{1}{2} [(A\mathbf{x}, \mathbf{x}) - 2(\mathbf{b}, \mathbf{x}) + (A\mathbf{x}^*, \mathbf{x}^*)] \\ &\quad - (A\mathbf{x}^* - \mathbf{b}, \mathbf{x}^*) \\ &= -\frac{1}{2} [(A\mathbf{x}, \mathbf{x}) - 2(A\mathbf{x}^*, \mathbf{x}) + (A\mathbf{x}^*, \mathbf{x}^*)] \\ &= -\frac{1}{2} (A(\mathbf{x} - \mathbf{x}^*), \mathbf{x} - \mathbf{x}^*) \geqslant 0, \quad (7.5.4) \end{aligned}$$

即 \mathbf{x}^* 是 $F(\mathbf{x})$ 的唯一极小点。

定理表明, 求线性方程组 (7.5.1) 的解的问题, 可以转化成求函数 $F(\mathbf{x})$ 的极小点。求 $F(\mathbf{x})$ 的极小点的计算过程, 大体可以归结为: 从任意向量 \mathbf{x} 出发, 沿着某一合适的方向向量 \mathbf{p} 进行修改, 从而得到一个新的近似解向量 $\bar{\mathbf{x}}$, 使得 $F(\mathbf{x}) > F(\bar{\mathbf{x}})$, 如此不断地修改下去, 最终得到 $F(\mathbf{x})$ 的极小点 \mathbf{x}^* 。从计算的过程可以看出, 求 $F(\mathbf{x})$ 的极小点 \mathbf{x}^* 的问题, 关键是选择修正方向向量。

5.2 最速下降法

最速下降法是一种逐步逼近的方法, 每次的逼近值都是从 $F(\mathbf{x})$ 的值的减小最快的方向上选取, 由此而命名这个方法。其具体迭代公式推导如下:

首先任取一个向量 $\mathbf{x}^{(0)}$ 作为初始向量, 选取一个方向 \mathbf{p}_0 使 $F(\mathbf{x})$ 的值沿方向 \mathbf{p}_0 减小最快, 从 (7.5.3) 知道

$$\text{grad} F(\mathbf{x}^{(0)}) = A\mathbf{x}^{(0)} - \mathbf{b} = -\mathbf{r}^{(0)},$$

即 $F(\mathbf{x})$ 从 $\mathbf{x}^{(0)}$ 出发沿 $\mathbf{r}^{(0)}$ 方向减小最快。所以选取 $\mathbf{p}_0 = \mathbf{r}^{(0)}$, 再在 $\mathbf{x}^{(0)} + a\mathbf{p}^{(0)} = \mathbf{x}^{(0)} + a\mathbf{r}^{(0)}$ 上选取一点 $\mathbf{x}^{(1)}$ 使

$$F(\mathbf{x}^{(1)}) = \min_a F(\mathbf{x}^{(0)} + a\mathbf{r}^{(0)}).$$

设

$$\begin{aligned} \phi(a) &= F(\mathbf{x}^{(0)} + a\mathbf{r}^{(0)}) \\ &= \frac{1}{2} (A(\mathbf{x}^{(0)} + a\mathbf{r}^{(0)}), (\mathbf{x}^{(0)} + a\mathbf{r}^{(0)})) \\ &\quad - (\mathbf{b}, \mathbf{x}^{(0)} + a\mathbf{r}^{(0)}) \\ &= \frac{1}{2} a^2 (A\mathbf{r}^{(0)}, \mathbf{r}^{(0)}) - a (\mathbf{r}^{(0)}, \mathbf{r}^{(0)}) \\ &\quad + F(\mathbf{x}^{(0)}). \end{aligned}$$

从方程

$$\frac{d\phi(\alpha)}{d\alpha} = \alpha (Ar^{(0)}, r^{(0)}) - (r^{(0)}, r^{(0)}) = 0$$

中解得

$$\alpha = (r^{(0)}, r^{(0)}) / (Ar^{(0)}, r^{(0)}) = \alpha_0,$$

又因为 $\frac{d^2\phi(\alpha)}{d\alpha^2} = (Ar^{(0)}, r^{(0)}) > 0$, 所以

$$\min_{\alpha} F(x^{(0)} + \alpha r^{(0)}) = F(x^{(0)} + \alpha_0 r^{(0)}).$$

于是解的第一次近似为

$$\begin{cases} x^{(1)} = x^{(0)} + \alpha_0 r^{(0)}, \\ \alpha_0 = (r^{(0)}, r^{(0)}) / (Ar^{(0)}, r^{(0)}), \\ r^{(0)} = b - Ax^{(0)}. \end{cases} \quad (7.5.5)$$

重复上述过程, 如果第 k 次近似向量 $x^{(k)}$ 已经求得, 那末, 从 $x^{(k)}$ 出发, 则使 $F(x)$ 之值减少最快的方向为 $r^{(k)} = b - Ax^{(k)}$, 同样可以求出

$$\min_{\alpha} F(x^{(k)} + \alpha r^{(k)}) = F(x^{(k)} + \alpha_k r^{(k)}),$$

其中

$$\alpha_k = (r^{(k)}, r^{(k)}) / (Ar^{(k)}, r^{(k)}).$$

于是第 $k+1$ 次近似为

$$\begin{cases} x^{(k+1)} = x^{(k)} + \alpha_k r^{(k)}, \\ \alpha_k = \frac{(r^{(k)}, r^{(k)})}{(Ar^{(k)}, r^{(k)})}, \quad r^{(k)} = b - Ax^{(k)}, \end{cases} \quad (7.5.6)$$

而且满足如下关系式:

$$\begin{aligned} F(x^{(k+1)}) - F(x^{(k)}) &= F(x^{(k)} + \alpha_k r^{(k)}) - F(x^{(k)}) \\ &= \frac{1}{2} \alpha_k^2 (Ar^{(k)}, r^{(k)}) - \alpha_k (r^{(k)}, r^{(k)}) \\ &= \frac{1}{2} - \frac{(r^{(k)}, r^{(k)})^2}{(Ar^{(k)}, r^{(k)})} = -\frac{(r^{(k)}, r^{(k)})^2}{(Ar^{(k)}, r^{(k)})} \end{aligned}$$

$$= -\frac{1}{2} \frac{(r^{(k)}, r^{(k)})^2}{(Ar^{(k)}, r^{(k)})} < 0, \quad (7.5.7)$$

即

$$F(x^{(k+1)}) < F(x^{(k)}).$$

这个结果说明了数列 $\{F(x^{(k)})\}$ 是单调下降的。但是，是否有

$$F(x^{(k)}) \rightarrow F(x^*),$$

还需要加以证明。为此先证明下列引理。

引理 7.5.1 设正数 λ_n, λ_1 分别为实对称正定阵 A 的最大与最小特征值，则不等式：

$$F(x^{(k+1)}) - F(x^*) \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^{2(k+1)} [F(x^{(k)}) - F(x^*)] \quad (7.5.8)$$

成立。

证明 根据(7.5.4)和 $Ax^* = b$ 可以证明

$$\begin{aligned} F(x^{(k)}) - F(x^*) &= -\frac{1}{2} (A(x^{(k)} - x^*), x^{(k)} - x^*) \\ &= -\frac{1}{2} (Ax^{(k)} - b, A^{-1}(Ax^{(k)} - b)) \\ &= -\frac{1}{2} (r^{(k)}, A^{-1}r^{(k)}), \end{aligned} \quad (7.5.9)$$

又由(7.5.7)，可得

$$\frac{F(x^{(k)}) - F(x^*)}{F(x^{(k)}) - F(x^{(k+1)})} = \frac{(r^{(k)}, A^{-1}r^{(k)}) (Ar^{(k)}, r^{(k)})}{(r^{(k)}, r^{(k)})^2}. \quad (7.5.10)$$

根据线性代数的知识，当 A 为对称正定阵时， A 的特征值可以排成： $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ ，而且一定存在一组数 $\beta_1, \beta_2, \dots, \beta_n$ 使(7.5.10)化为如下形式：

$$\begin{aligned}
\frac{(r^{(k)}, A^{-1}r^{(k)}) (Ar^{(k)}, r^{(k)})}{(r^{(k)}, r^{(k)})^2} &= \frac{\sum \beta_i^2 \lambda_i^{-1} \sum \beta_i^2 \lambda_i}{(\sum \beta_i^2)^2} \\
&= \frac{1}{(\sum \beta_i^2)^2} \sum \beta_i^2 \sqrt{\frac{\lambda_n \lambda_1}{\lambda_i}} \cdot \sum \beta_i^2 \frac{\lambda_i}{\sqrt{\lambda_n \lambda_1}} \\
&\leq \frac{1}{4(\sum \beta_i^2)^2} \left[\sum \beta_i^2 \left(\sqrt{\frac{\lambda_n \lambda_1}{\lambda_i}} + \frac{\lambda_i}{\sqrt{\lambda_n \lambda_1}} \right) \right]^2. \quad (7.5.11)
\end{aligned}$$

设 $\sqrt{\frac{\lambda_n \lambda_1}{\lambda_i}} = y$, 则 $\sqrt{\frac{\lambda_1}{\lambda_n}} \leq y \leq \sqrt{\frac{\lambda_n}{\lambda_1}}$. 因为 $y + \frac{1}{y}$ 在 $y \in \left[\sqrt{\frac{\lambda_1}{\lambda_n}}, \sqrt{\frac{\lambda_n}{\lambda_1}} \right]$ 的最大值只能在 y 的区间端点 $\sqrt{\frac{\lambda_1}{\lambda_n}}$ 或 $\sqrt{\frac{\lambda_n}{\lambda_1}}$ 上取得, 而在两端点上取值均为 $\sqrt{\frac{\lambda_1}{\lambda_n}} + \sqrt{\frac{\lambda_n}{\lambda_1}}$, 从而 (7.5.11) 可以化成:

$$\begin{aligned}
\frac{(r^{(k)}, A^{-1}r^{(k)}) (Ar^{(k)}, r^{(k)})}{(r^{(k)}, r^{(k)})^2} &\leq \frac{1}{4(\sum \beta_i^2)^2} \left[\sum \beta_i^2 \left(\sqrt{\frac{\lambda_1}{\lambda_n}} \right. \right. \\
&\quad \left. \left. + \sqrt{\frac{\lambda_n}{\lambda_1}} \right) \right]^2 = \frac{(\lambda_1 + \lambda_n)^2}{4\lambda_1 \lambda_n} = \frac{1}{e}. \quad (7.5.12)
\end{aligned}$$

由 (7.5.10) 和 (7.5.12) 得到

$$F(x^{(k)}) - F(x^{(k+1)}) \geq e(F(x^{(k)}) - F(x^*)).$$

当 $\lambda_1 \neq \lambda_n$ 时, $0 < e < 1$, 于是

$$\begin{aligned}
F(x^{(k+1)}) - F(x^*) &= [F(x^{(k)}) - F(x^*)] - [F(x^{(k)}) - F(x^{(k+1)})] \\
&\leq [F(x^{(k)}) - F(x^*)] - e[F(x^{(k)}) - F(x^*)] \\
&= (1-e)[F(x^{(k)}) - F(x^*)]. \quad (7.5.13)
\end{aligned}$$

重复使用不等式 (7.5.13), 最后可得

$$F(x^{(k+1)}) - F(x^*) \leq (1-e)^{k+1} [F(x^{(0)}) - F(x^*)]$$

$$\begin{aligned}
&= \left(1 - \frac{4\lambda_1\lambda_n}{(\lambda_1 + \lambda_n)^2}\right)^{k+1} [F(\mathbf{x}^{(0)}) - F(\mathbf{x}^*)] \\
&= \left(\frac{\lambda_n - \lambda_1}{\lambda_1 + \lambda_n}\right)^{2(k+1)} [F(\mathbf{x}^{(0)}) - F(\mathbf{x}^*)],
\end{aligned}$$

于是(7.5.8)得证.

定理 7.5.2 设 A 为对称正定阵, 则

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^{k+1} \sqrt{\frac{2l}{\lambda_1}}. \quad (7.5.14)$$

其中 $l = F(\mathbf{x}^{(0)}) - F(\mathbf{x}^*)$.

证明 因为 $(A\mathbf{x}, \mathbf{x}) = \sum \beta_i^2 \lambda_i \geq \sum \lambda_1 \beta_i^2 = \lambda_1 (\mathbf{x}, \mathbf{x})$,
所以

$$(\mathbf{x}, \mathbf{x}) \leq \frac{1}{\lambda_1} (A\mathbf{x}, \mathbf{x}),$$

于是

$$\begin{aligned}
\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 &= (\mathbf{x}^{(k+1)} - \mathbf{x}^*, \mathbf{x}^{(k+1)} - \mathbf{x}^*) \\
&\leq \frac{1}{\lambda_1} (A(\mathbf{x}^{(k+1)} - \mathbf{x}^*), \mathbf{x}^{(k+1)} - \mathbf{x}^*) \\
&= \frac{2}{\lambda_1} (F(\mathbf{x}^{(k+1)}) - F(\mathbf{x}^*)).
\end{aligned}$$

再根据引理7.5.1的结果, 得到

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 \leq \frac{2}{\lambda_1} \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^{2(k+1)} (F(\mathbf{x}^{(0)}) - F(\mathbf{x}^*)),$$

即

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^{k+1} \sqrt{\frac{2l}{\lambda_1}}.$$

因为 $\lambda_n \geq \lambda_1 > 0$, 所以 $0 \leq \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} < 1$, 于是当 $k \rightarrow \infty$ 时, $\mathbf{x}^{(k+1)} \rightarrow \mathbf{x}^*$, 即最速下降法是收敛的。但是, 如果 λ_n 与 λ_1 相差很大, 即 $(\lambda_n - \lambda_1)/(\lambda_n + \lambda_1) \approx 1$ 时, 由 (7.5.14) 知道, 最速下降法收敛得很慢。实际计算也表明最速下降法一般收敛较慢, 特别对于病态方程组更是如此。因此, 必须设法提高最速下降法的收敛速度。

5.3 共轭斜量法

解线性方程组的最速下降法, 由于它具有局部性质, 即在 \mathbf{x} 附近函数 $F(\mathbf{x})$ 沿方向 \mathbf{r} 下降较快, 但是从总体来说并不如此, 因此收敛速度并不理想, 目前使用得较少。下面我们将选择一些方向, 使得 $F(\mathbf{x})$ 沿着这些方向能更快地逼近它的最小值。由于选择的方向具有“共轭”的性质, 所以称为共轭斜量法。

在构造共轭斜量法的计算公式之前, 先引入关于 A 一共轭 (正交) 的概念。

设 A 为对称正定矩阵, \mathbf{x} 、 \mathbf{y} 是实的列向量, 如果

$$(\mathbf{x}, A\mathbf{y}) = (A\mathbf{x}, \mathbf{y}) = 0, \quad (7.5.15)$$

则称向量 \mathbf{x} 和 \mathbf{y} 是 A 一共轭 (正交) 的。

下面推导共轭斜量法的计算公式。

对于任意的初始向量 $\mathbf{x}^{(0)}$, 有剩余向量

$$\mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)}.$$

从 $\mathbf{x}^{(0)}$ 出发沿 $\mathbf{p}^{(0)} = \mathbf{r}^{(0)}$ 寻找新的近似向量 $\mathbf{x}^{(1)}$, 取

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{r}^{(0)},$$

$$\alpha_0 = (\mathbf{r}^{(0)}, \mathbf{p}^{(0)}) / (\mathbf{p}^{(0)}, A\mathbf{p}^{(0)}),$$

再从 $\mathbf{x}^{(1)}$ 出发沿 $\mathbf{p}^{(1)} = \mathbf{r}^{(1)} + \beta_0 \mathbf{p}^{(0)}$ (其中 $\mathbf{r}^{(1)} = \mathbf{b} - A\mathbf{x}^{(1)}$) 寻找新的近似向量 $\mathbf{x}^{(2)}$, 并要求 $\mathbf{p}^{(1)}$ 与 $\mathbf{p}^{(0)}$ 成为 A 一共轭, 从而

$$\beta_1 = -(\mathbf{r}^{(1)}, \mathbf{A}\mathbf{p}^{(0)})/(\mathbf{p}^{(0)}, \mathbf{A}\mathbf{p}^{(0)}).$$

于是 $\mathbf{x}^{(2)}$ 可取为

$$\begin{cases} \mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{p}^{(1)}, \\ \alpha_1 = (\mathbf{r}^{(1)}, \mathbf{p}^{(1)})/(\mathbf{p}^{(1)}, \mathbf{A}\mathbf{p}^{(1)}). \end{cases}$$

如果 $\mathbf{x}^{(i)}$ 已经求得, 再去寻找新的近似解向量时, 可以从 $\mathbf{x}^{(i)}$ 出发选取方向向量 $\mathbf{p}^{(i)} = \mathbf{r}^{(i)} + \beta_{i-1} \mathbf{p}^{(i-1)}$, 并要求 $\mathbf{p}^{(i)}$ 与 $\mathbf{p}^{(i-1)}$ 成为 \mathbf{A} —共轭, 即

$$(\mathbf{A}\mathbf{p}^{(i)}, \mathbf{p}^{(i-1)}) = (\mathbf{p}^{(i)}, \mathbf{A}\mathbf{p}^{(i-1)}) = 0,$$

从而

$$\beta_{i-1} = -(\mathbf{r}^{(i)}, \mathbf{A}\mathbf{p}^{(i-1)})/(\mathbf{p}^{(i-1)}, \mathbf{A}\mathbf{p}^{(i-1)}), \quad (7.5.16)$$

其中 $\mathbf{r}^{(i)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(i)}$, 于是新的近似解向量 $\mathbf{x}^{(i+1)}$ 可取为

$$\begin{cases} \mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} + \alpha_i \mathbf{p}^{(i)}, \\ \alpha_i = (\mathbf{r}^{(i)}, \mathbf{p}^{(i)})/(\mathbf{p}^{(i)}, \mathbf{A}\mathbf{p}^{(i)}). \end{cases} \quad (7.5.17)$$

以上就是共轭斜量法的主要计算步骤, 实际计算中, 我们还可以将上面的公式进行简化. 为此, 首先说明向量系 $\{\mathbf{r}^{(i)}\}$ 与 $\{\mathbf{p}^{(i)}\}$ 之间的关系和构造方法. 因为

$$\begin{cases} \mathbf{r}^{(i+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(i+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(i)} - \alpha_i \mathbf{A}\mathbf{p}^{(i)} = \mathbf{r}^{(i)} - \alpha_i \mathbf{A}\mathbf{p}^{(i)}, \\ \mathbf{p}^{(i+1)} = \mathbf{r}^{(i+1)} + \beta_i \mathbf{p}^{(i)}, \end{cases} \quad (7.5.18)$$

所以当 $i = 0$ 时,

$$\begin{cases} \mathbf{r}^{(1)} = \mathbf{r}^{(0)} - \alpha_0 \mathbf{A}\mathbf{p}^{(0)}, \\ \mathbf{p}^{(1)} = \mathbf{r}^{(1)} + \beta_0 \mathbf{p}^{(0)} = \mathbf{r}^{(1)} + \beta_0 \mathbf{r}^{(0)} \end{cases}$$

当 $i = 1$ 时,

$$\begin{cases} \mathbf{r}^{(2)} = \mathbf{r}^{(1)} - \alpha_1 \mathbf{A}\mathbf{p}^{(1)} = \mathbf{r}^{(0)} - \alpha_0 \mathbf{A}\mathbf{p}^{(0)} - \alpha_1 \mathbf{A}\mathbf{p}^{(1)}, \\ \mathbf{p}^{(2)} = \mathbf{r}^{(2)} + \beta_1 \mathbf{p}^{(1)} = \mathbf{r}^{(2)} + \beta_1 \mathbf{r}^{(1)} + \beta_1 \beta_0 \mathbf{r}^{(0)}. \end{cases}$$

继续上述推导可得

$$r^{(i+1)} = r^{(0)} - \alpha_0 A p^{(0)} - \alpha_1 A p^{(1)} - \dots - \alpha_i A p^{(i)}, \quad (7.5.19)$$

$$p^{(i+1)} = r^{(i+1)} + \beta_1 r^{(i)} + \beta_i \beta_{i-1} r^{(i-1)} + \dots + \beta_i \beta_{i-1} \dots \beta_0 r^{(0)}, \quad (7.5.20)$$

$$r^{(i+1)} = p^{(i+1)} - \beta_i p^{(i)}. \quad (7.5.21)$$

有了关系式(7.5.19), (7.5.20), (7.5.21)后, 我们可以证明如下定理.

定理 7.5.3 剩余向量系 $\{r^{(i)}\}$ 构成一个正交向量组, 即

$$(r^{(i)}, r^{(j)}) = 0 \quad (i \neq j), \quad (7.5.22)$$

方向向量系 $\{p^{(i)}\}$ 构成一个 A —共轭向量系, 即

$$(p^{(i)}, A p^{(j)}) = 0 \quad (i \neq j). \quad (7.5.23)$$

证明 应用归纳法来证明.

$$(r^{(0)}, r^{(1)}) = (r^{(0)}, r^{(0)} - \alpha_0 A p^{(0)}) = (r^{(0)}, r^{(0)}) - (r^{(0)}, \alpha_0 A p^{(0)})$$

$$= (r^{(0)}, r^{(0)}) - \frac{(r^{(0)}, r^{(0)})}{(r^{(0)}, A r^{(0)})} (r^{(0)}, A r^{(0)}) = 0,$$

$$(p^{(0)}, A p^{(1)}) = (r^{(0)}, A(r^{(1)} + \beta_0 r^{(0)})) = (r^{(0)}, A r^{(1)})$$

$$+ \beta_0 (r^{(0)}, A r^{(0)}) = (r^{(0)}, A r^{(1)}) - \frac{(r^{(1)}, A r^{(0)})}{(r^{(0)}, A r^{(0)})} (r^{(0)}, A r^{(0)})$$

$$= (r^{(0)}, A r^{(1)}) - (r^{(1)}, A r^{(0)}) = 0.$$

设 $r^{(0)}, r^{(1)}, \dots, r^{(k)}$ 已两两正交; $p^{(0)}, p^{(1)}, \dots, p^{(k)}$ 已两两 A —共轭, 现在要证

$$(r^{(k+1)}, r^{(j)}) = 0 \quad (j = 0, 1, \dots, k).$$

因为

$$\begin{aligned} (r^{(k+1)}, r^{(j)}) &= (r^{(k)} - \alpha_k A p^{(k)}, r^{(j)}) \\ &= (r^{(k)}, r^{(j)}) - \alpha_k (A p^{(k)}, r^{(j)}) \\ &= (r^{(k)}, r^{(j)}) - \alpha_k (A p^{(k)}, p^{(j)} - \beta_{j-1} p^{(j-1)}) \\ &= (r^{(k)}, r^{(j)}) - \alpha_k (A p^{(k)}, p^{(j)}) \\ &\quad + \alpha_k \beta_{j-1} (A p^{(k)}, p^{(j-1)}) = 0, \end{aligned}$$

所以(7.5.22)成立。再证

$$(p^{(k+1)}, Ap^{(j)}) = 0 \quad (j=0, 1, \dots, k).$$

因为

$$\begin{aligned} (p^{(k+1)}, Ap^{(j)}) &= (r^{(k+1)} + \beta_k p^{(k)}, Ap^{(j)}) \\ &= (r^{(k+1)}, Ap^{(j)}) + \beta_k (p^{(k)}, Ap^{(j)}), \end{aligned}$$

所以当 $j=k$ 时,

$$\begin{aligned} (p^{(k+1)}, Ap^{(k)}) &= (r^{(k+1)}, Ap^{(k)}) \\ &\quad - \frac{(r^{(k+1)}, Ap^{(k)})}{(p^{(k)}, Ap^{(k)})} \cdot (p^{(k)}, Ap^{(k)}) = 0. \end{aligned}$$

当 $j \leq k-1$ 时,

$$\begin{aligned} (p^{(k+1)}, Ap^{(j)}) &= (r^{(k+1)}, Ap^{(j)}) \\ &= (r^{(k+1)}, \frac{1}{a_j} (r^{(j)} - r^{(j+1)})) \\ &= \frac{1}{a_j} [(r^{(k+1)}, r^{(j)}) - (r^{(k+1)}, r^{(j+1)})] \\ &= 0, \end{aligned}$$

从而证明了(7.5.23)。

下面将 a_i 和 β_i 的计算公式进行简化。由(7.5.20)和(7.5.22)可得

$$\begin{aligned} (r^{(i)}, p^{(i)}) &= (r^{(i)}, r^{(i)} + \beta_{i-1} r^{(i-1)} + \dots + \beta_{i-1} \dots \beta_0 r^{(0)}) \\ &= (r^{(i)}, r^{(i)}), \end{aligned}$$

于是,

$$a_i = (r^{(i)}, p^{(i)}) / (p^{(i)}, Ap^{(i)}) = (r^{(i)}, r^{(i)}) / (p^{(i)}, Ap^{(i)}),$$

所以当 $r^{(i)} \neq 0$ 时, $a_i > 0$ 。又由(7.5.18), (7.5.22), (7.5.23)可得

$$\begin{aligned} (r^{(i)}, Ap^{(i-1)}) &= (r^{(i)}, \frac{1}{a_{i-1}} (r^{(i-1)} - r^{(i)})) \\ &= - (r^{(i)}, r^{(i)}) / a_{i-1}, \end{aligned}$$

$$\begin{aligned}
(p^{(i-1)}, Ap^{(i-1)}) &= (r^{(i-1)} + \beta_{i-2} p^{(i-2)}, Ap^{(i-1)}) \\
&= (r^{(i-1)}, Ap^{(i-1)}) \\
&= (r^{(i-1)}, \frac{1}{a_{i-1}} (r^{(i-1)} - r^{(i)})) \\
&= (r^{(i-1)}, r^{(i-1)}) / a_{i-1},
\end{aligned}$$

于是

$$\beta_{i-1} = - \frac{(r^{(i)}, Ap^{(i-1)})}{(p^{(i-1)}, Ap^{(i-1)})} = \frac{(r^{(i)}, r^{(i)})}{(r^{(i-1)}, r^{(i-1)})}.$$

可见当 $r^{(i)} \neq 0$ 时, $\beta_{i-1} > 0$, 从而得到共轭斜量法的算法.

算法 7.5.1 设 A 为 n 阶对称正定阵, $b = (b_1, b_2, \dots, b_n)^T$, eps 为 (r, r) 或 (Ap, p) 的允许误差, 引入比例因子 d 以防发生溢出; 取初始向量为 0 和限定迭代次数不超过 $3n$.

1) $d \leftarrow 2^\omega \triangleq \text{Max}(|b_i|)$, 其中 ω 为整数.

2) $p = r = b/d$.

3) 对 $k = 0, 1, \dots, 3n$,

3.1) $q \leftarrow d \cdot (r, r) / (Ap, p)$;

3.2) 对 $i = 1, 2, \dots, n$,

$$x_i \leftarrow x_i + qp_i,$$

$$\bar{r}_i \leftarrow r_i - q \sum_{j=1}^n a_{ij} p_j, \quad (r, r) \leq \text{eps 转 6),}$$

3.3) $e \leftarrow (\bar{r}, \bar{r}) / (r, r), \quad r \leftarrow \bar{r}$,

3.4) $p \leftarrow r + ep, (Ap, p) \leq \text{eps 转 6),}$

3.5) NEXT k .

4) $\bar{d} \leftarrow 2^\omega \triangleq \max(|p_i|)$.

5) $p \leftarrow p/\bar{d}, r \leftarrow r/\bar{d}, d \leftarrow \bar{d}d$ 转 3).

6) 输出结果 x 和 k .

例 7.5.1 用共轭斜量法解方程组

$$\begin{cases} 3x_1 + x_2 = 5; \\ x_1 + 2x_2 = 5. \end{cases}$$

解 $A = \begin{bmatrix} 3 & 1 \\ 1 & 2 \end{bmatrix}, \quad b = \begin{bmatrix} 5 \\ 5 \end{bmatrix}.$

取 $x^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix},$

1) 计算 $r^{(0)} = b - Ax^{(0)} = p^{(0)} = \begin{bmatrix} 5 \\ 5 \end{bmatrix},$

$$Ap^{(0)} = \begin{bmatrix} 20 \\ 15 \end{bmatrix}, \quad (p^{(0)}, Ap^{(0)}) = 175,$$

$$(r^{(0)}, r^{(0)}) = 50, \quad \alpha_0 = 2/7.$$

得 $x^{(1)} = x^{(0)} + \alpha_0 p^{(0)} = \frac{2}{7} \begin{bmatrix} 5 \\ 5 \end{bmatrix}.$

2) $r^{(1)} = \begin{bmatrix} -5 \\ 1 \end{bmatrix}, \quad \beta_0 = 1/49,$

$$p^{(1)} = r^{(1)} + \beta_0 p^{(0)} = \frac{1}{49} \begin{bmatrix} -30 \\ 40 \end{bmatrix},$$

$$Ap^{(1)} = \frac{1}{49} \begin{bmatrix} -50 \\ 50 \end{bmatrix}, \quad (r^{(1)}, r^{(1)}) = \frac{50}{49}.$$

$$(p^{(1)}, Ap^{(1)}) = \frac{3500}{49^2}, \quad \alpha_1 = \frac{7}{10},$$

得 $x^{(2)} = x^{(1)} + \alpha_1 p^{(1)} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} = x^*,$

这时 $r^{(2)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$

下面我们将证明, 用共轭斜量法所求得 $x^{(i)}$, 随着 i 的增大而一步步地逼近方程组的精确解, 而且若在计算中没有误差

时, 至多到 n 步便能得到方程组的精确解。

定理 7.5.4 依共轭斜量法求解 n 阶线性方程组, 最多计算 n 步便可得到方程组 (7.5.1) 的精确解。

证明 根据定理 7.5.3 可知, 如果 $n+1$ 个剩余向量 $r^{(0)}, r^{(1)}, \dots, r^{(n)}$ 均不为零, 那末这些向量必线性无关, 从而与线性无关向量组的性质发生矛盾, 故在 $r^{(0)}, r^{(1)}, \dots, r^{(n)}$ 中至少有一个 $r^{(k)}$ 为零, 由 $r^{(k)} = b - Ax^{(k)} = 0$, 得 $x^{(k)} = x^*$ 。

定理 7.5.5 设当 $l < k$ 时, $x^{(k)}$ 比 $x^{(l)}$ 更接近于精确解 x^* 。即

$$\|x^{(l)} - x^*\| \geq \|x^{(k)} - x^*\|。$$

证明 如果能证明

$$\|x^{(k)} - x^*\| \geq \|x^{(k+1)} - x^*\|, \quad (7.5.24)$$

那末, 用不等式进行递推, 便能得到定理 7.5.5 的结论。

因为 $x^{(k+1)} = x^{(k)} + \alpha_k p^{(k)}$, 所以

$$\begin{aligned} \|x^{(k)} - x^*\|^2 &= (x^{(k+1)} - x^* - \alpha_k p^{(k)}, x^{(k+1)} - x^* - \alpha_k p^{(k)}) \\ &= \|x^{(k+1)} - x^*\|^2 - 2\alpha_k (p^{(k)}, x^{(k+1)} - x^*) + \alpha_k^2 (p^{(k)}, p^{(k)})。 \end{aligned} \quad (7.5.25)$$

由于 $\alpha_k \geq 0$, $(p^{(k)}, p^{(k)}) > 0$, 所以只要证明 $(p^{(k)}, x^{(k+1)} - x^*) \leq 0$, 定理的结论就成立了。因为

$$\begin{aligned} x^{(k+1)} - x^* &= x^{(k)} + \alpha_k p^{(k)} - x^* = x^{(k-1)} + \alpha_{k-1} p^{(k-1)} + \alpha_k p^{(k)} - x^* = \dots \\ &= x^{(0)} + \alpha_0 p^{(0)} + \alpha_1 p^{(1)} + \dots + \alpha_k p^{(k)}, \end{aligned}$$

据定理 7.5.4 必有 $x^{(m)} = x^*$ ($m \leq n$), 所以 $x^{(m)} = x^* = x^{(0)} + \alpha_0 p^{(0)} + \dots + \alpha_{m-1} p^{(m-1)}$ 。因为当 $k = m-1$ 时, (7.5.24) 显然成立。于是只需讨论 $k < m-1$ 的情况。此时, 由

$$\begin{aligned} x^{(k+1)} - x^* &= -\alpha_{k+1} p^{(k+1)} - \alpha_{k+2} p^{(k+2)} - \dots - \alpha_{m-1} p^{(m-1)}, \\ \text{则 } (p^{(k)}, x^{(k+1)} - x^*) &= -\alpha_{k+1} (p^{(k)}, p^{(k+1)}) - \alpha_{k+2} (p^{(k)}, p^{(k+2)}) \\ &\quad - \dots - \alpha_{m-1} (p^{(k)}, p^{(m-1)}). \end{aligned}$$

又由于 $a_{h+1} \geq 0$ ，而且容易证明

$$(p^{(k)}, p^{(k+1)}) = \|r^{(k)}\|^2 \|r^{(k+1)}\|^2 \sum_{j=0}^k (1/\|r^{(j)}\|^2) \geq 0, \quad (7.5.26)$$

从而 $(p^{(k)}, x^{(k+1)} - x^*) \leq 0$ ，即本定理成立。

共轭斜量法虽然本质上是一种直接法，在有限步内可以得到方程组的精确解。可是，在实际计算中由于存在舍入误差，使剩余向量不能精确满足正交关系，所以可能有 $r^{(n)} \neq 0$ ，而且当系数矩阵病态时， $r^{(n)}$ 偏离零的程度也更厉害。不过，由于共轭斜量法的计算公式具有迭代法的特点，如果 $r^{(n)} \neq 0$ ，还可以用此法继续计算下去，只要 $r^{(n+1)} \neq 0$ ， $F(x)$ 的值总会减小，直到由于舍入误差的影响使近似解不能再改进为止。因此，最后总可以得到更好一些的近似解。

本章讨论了解线性方程组的线性与非线性迭代法，重点介绍了 SOR 法，Chebyshev 半迭代法和共轭斜量法，这三种方法也是常用的方法。

Jacobi 法与 Gauss-Siedel 法的收敛性并不互相包含，不过当两者皆收敛时，往往后者收敛快。当收敛慢时，可用 Chebyshev 加速，从而得到 Chebyshev 半迭代法。当原迭代矩阵 B 的特征值全部是实数时，Chebyshev 加速是 very 有效的。但是当原方程的系数矩阵 A 是具有相容次序的对称正定阵时，SOR 方法的收敛速度比简单迭代法的 Chebyshev 加速还要快两倍。即使矩阵不具有相容次序，在某些条件下也还比 Jacobi 方法或 Gauss-Siedel 方法快一个数量级，因此，在实用中常常使用 SOR 方法，其松弛因子的选择往往依赖于实际经验。

当系数矩阵 A 是对称正定阵时，共轭斜量法是一种较好的方法。它不需要估计参数，但是存在着计算过程中舍入误差比较敏感的缺点。近年来提出的“不完全分解预处理共轭斜量

法”，弥补了这方面的缺点。这种方法虽然在理论上还不成熟，但是在某些实际应用中，却取得了惊人的效果，受到了国内外的普遍重视。

近年来，对迭代法还常常用于处理法的技巧来减小原矩阵的条件数以加快迭代法的收敛速度。尤其对病态方程组效果特别显著。另外，还发展了不少对某些具体问题的解决特别有效的迭代法，如“强隐式迭代法”等。最后还应指出，迭代法的使用与具体问题的特点有着很密切的联系。在各种实际问题中使用迭代法的经验与针对问题进行详细讨论是很重要的。

第七章 习 题

7.1 设线性方程组 $Ax = b$ 中， A 为

$$1) \begin{pmatrix} 3 & 1 & 2 \\ 0 & 2 & 1 \\ 1 & 0 & 1 \end{pmatrix}, \quad 2) \begin{pmatrix} 1 & 2 & -2 \\ 1 & 1 & 1 \\ 2 & 2 & 1 \end{pmatrix},$$

若用 Jacobi 迭代法，Gauss-Siedel 迭代法求解，是否收敛？若收敛，求方程组的解。

7.2 设 $Ax = b$ ，其中 $A \in R^{n \times n}$ ，试证明

$$\max_i \sum_{j \neq i} \frac{|a_{ji}|}{|a_{ii}|} < 1$$

是 Jacobi 迭代法收敛的充分条件。

7.3 用逐次超松弛迭代法解方程组（取 $\omega = 0.9$ ）

$$\begin{cases} 5x_1 + 2x_2 + x_3 = -12, \\ -x_1 + 4x_2 + 2x_3 = 20, \\ 2x_1 - 3x_2 + 10x_3 = 3 \end{cases}$$

准确到小数点后4位。

7.4 证明定理7.2.3和定理7.2.5的充分性。

7.5 试证明 n 阶矩阵 A 不可约的充要条件是

$n=1$ 或者对任意两个不同的整数 i, j , 其中 $1 \leq i \leq n, 1 \leq j \leq n$, 则 $a_{ij} \neq 0$ 或者存在 i_1, i_2, \dots, i_s 满足

$$a_{ii_1} \cdot a_{i_1 i_2} \cdot \dots \cdot a_{i_{s-1} i_s} \neq 0.$$

7.6 证明 $\lim_{m \rightarrow \infty} (\|B^m\|)^{1/m} = S(B)$.

7.7 假若 A 为正定矩阵, 那末使用松弛因子 $\omega_1, \omega_2, \dots$ 的 SOR 方法收敛, 下列条件至少一个成立

1) 对某一 $\varepsilon > 0$, 对足够大的 i , 有

$$\varepsilon \leq \omega_i \leq 2 - \varepsilon \quad i = 1, 2, \dots.$$

2) 对所有足够大的 i , $0 \leq \omega_i \leq 2$ 和级数

$$\sum_{i=1}^{\infty} \omega_i (2 - \omega_i)$$

发散.

7.8 当 $\omega = 1$ 时, 对 (7.3.12) 使用简单迭代法, 证明该迭代法收敛.

7.9 设序列 $\{x_i\} = \{(-1)^{i+1}\}$, 证明序列 $\{x_i\}$ 为 Cesàro 可求和. 问是否存在 $\sum_{k=0}^n a_{nk} \quad n = 0, 1, 2, \dots$, 使作成的新序列 $y_n = \sum_{k=0}^n a_{nk} x_k$ 发散.

7.10 设方程组 $x = Bx + k$ 与 $Ax = b$ (其中 A 为非奇异矩阵) 等价, 且矩阵 B 的特征值 μ 满足: $-2 \leq \mu \leq 0.9, 1.2 \leq \mu \leq 3$. 使用线性定常迭代法

$$x^{(n+1)} = Bx^{(n)} + k,$$

构造一个收敛的半迭代法.

7.11 对习题 7.3 中的方程组分别利用基于 Jacobi 和 Gauss-Siedel 迭代法的 Chebyshev 迭代法求解.

7.12 证明公式 (7.5.26) 成立.

7.13 使用共轭斜量法解方程组

$$\begin{cases} 4x_1 - x_2 + 2x_3 = 12, \\ -x_1 + 5x_2 + 3x_3 = 10, \\ 2x_1 + 3x_2 + 6x_3 = 18. \end{cases}$$

7.14 设 $Ax = b$ 其中 A 为非奇异阵. 求证 $A^T A$ 为对称正定阵, 并写出用共轭斜量法解 $A^T A x = A^T b$ 的计算公式.

7.15 设 A 为对称正定阵, 求证 $Ax = b$ 的另一种共轭斜量法的迭代公式

为

$$\begin{cases} x^{(0)} \text{ 初始向量, } \Delta x_{-1} = \Delta r_{-1} = 0, \\ x^{(k+1)} = x^{(k)} + \Delta x_k, \\ \Delta x_k = \frac{1}{q_k} (r^{(k)} + e_k \Delta x_{k-1}), \\ q_k = \frac{(r^{(k)}, Ar^{(k)})}{(r^{(k)}, r^{(k)})} - e_k, \quad e_0 = 0, \\ e_k = \frac{(r^{(k)}, r^{(k)})}{(r^{(k-1)}, r^{(k-1)})} q_{k-1}, \\ r^{(k+1)} = r^{(k)} + \Delta r_k, \\ \Delta r_{k+1} = \frac{1}{q_k} (-Ar^{(k)} + e_k \Delta r_{k-1}), \\ k = 0, 1, 2, \dots \end{cases}$$

参 考 书

- [1] 徐树荣译, 线性代数方程组的计算机解法, 科学出版社.
- [2] D. M. Young. "Iterative Solution of large linear systems", Academic press, New York and London, 1971.
- [3] Stewart, G. "Introduction to Matrix Computation" Academic press. New York 1973.
- [4] Jennings, A. "Matrix computation for Engineers and Scientists", London New york. Sydney. Toronto, 1977.
- [5] 冯康等编, 数值计算方法, 国防工业出版社. 1978.
- [6] 曹志浩等编. 矩阵计算和方程求根. 1979.

第八章 非线性方程组的 牛顿迭代解法

§1 引 言

近年来，许多应用领域和数值分析本身都提出了多个实变量的非线性方程组

$$f_i(\xi_1, \xi_2, \dots, \xi_n) = 0 \quad (i = 1, \dots, n) \quad (8.1.1)$$

的数值求解问题。因此研究这一课题显得十分迫切。

非线性方程组 (8.1.1) 也可用向量形式表示：

$$F(x) = 0$$

这里 $F = (f_1, f_2, \dots, f_n)^T$, $x = (\xi_1, \xi_2, \dots, \xi_n)^T$. 向量函数 $F(x)$ 可视为一个从区域 $D \subset R^n$ 至 R^n 的映射，常记作 $F: D \subset R^n \rightarrow R^n$. 如果 f_1, f_2, \dots, f_n 皆为线性函数，则 (8.1.1) 即为 n 个变量的线性方程组；如果 $n=1$ ，则 (8.1.1) 即为单变量的非线性方程。这两个特例在数值分析教程中均作了较详细的研究。对于一般的非线性方程组的数值求解，很自然地会面临两个问题：

- (1) 方程组 (8.1.1) 的可解性分析；
- (2) 寻找有效的数值方法来求解 (8.1.1)。

所谓有效的意思，主要是指迭代算法具有较快的收敛速度和较少的运算量。第一个问题已远远超出本教程的范围，因为它属于非线性泛函分析的范畴。所以，以下我们均假定方程组 (8.1.1) 存在解。今后的任务仅在于构造各种迭代算法，以求得 (8.1.1) 的近似解。

对于许多迭代算法常成立以下定理：如果 x^* 为 (8.1.1) 在区域 D 中的一个解，而 x_0 为一个充分接近 $x^* \in D$ 的初始近似解，则由 x_0 出发产生的迭代序列 $\{x_k\}$ 恒有定义且收敛于 x^* 。这种类型的定理常称为局部收敛定理。固然，有极个别的情形，迭代序列 $\{x_k\}$ 收敛于 x^* 与初始近似解 x_0 的选取无关，这种收敛性质常称为整体收敛性。对于某迭代算法，若 n 维空间中存在球 $S(x^*, r) \subset D$ ，从该球中任何一点出发，由此生成的 $\{x_k\}$ 必收敛于 x^* ，则球 S 就称为该迭代算法的局部收敛域。非线性方程组的迭代解法的基本思想由以下三部分组成：

(1) 由一个极差的初始估猜 x_0 出发，设法进入局部收敛区域。在实际问题中，有时事先对所求的解已经有了一个较好的估计，则这一步就可以免去。

(2) 从局部收敛区域 S 中一个初始近似解 x_0 出发，通过加速收敛的迭代算法，求出足够精确的近似解。以后着重介绍的就是加速收敛的局部技巧，或者说局部方法。

(3) 继续求方程组的其余解。

§2 非线性方程组的牛顿迭代法

正如一维情形那样，牛顿法也是求解非线性方程组 (8.1.1) 的一个基本算法。许多其它的迭代算法都是在它的基础上作各种修改得到的。因此，无论在实践和理论方面它都具有重要意义。

一维牛顿法的几何意义是：在横坐标为 x_k 的点 P_k 引切线。于是在 P_k 点近傍以该切线近似替代曲线 $y = f(x)$ 。然后，将该切线与 x 轴的交点的横坐标 x_{k+1} 作为根 x^* 的新近似值。如此不断迭代，直至得到满意的近似值为止。现在，我们就来仿

照一维情形，构造出非线性方程组对应的牛顿迭代公式。

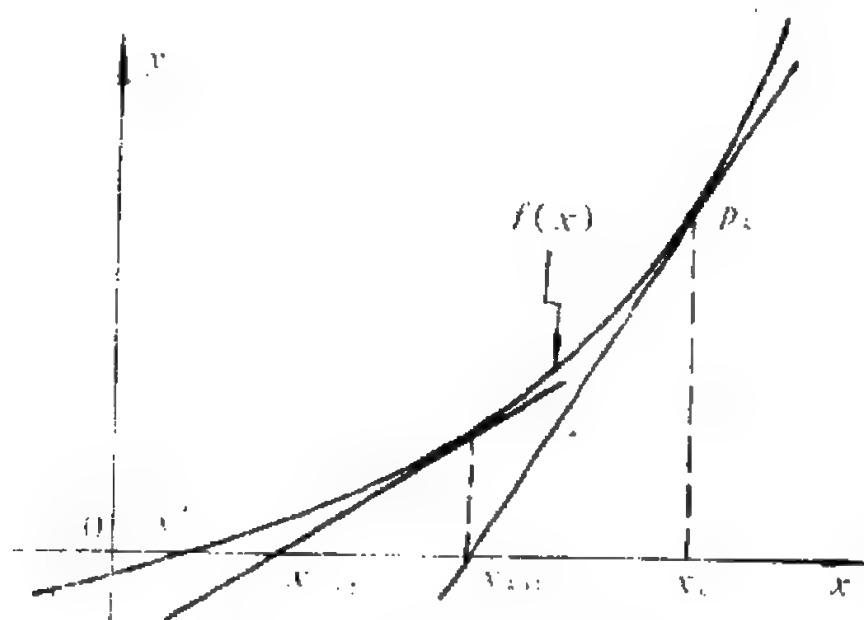


图8.1 牛顿法的几何意义

为叙述方便起见，仅对二个方程的方程组

$$\begin{cases} f_1(\xi_1, \xi_2) = 0, \\ f_2(\xi_1, \xi_2) = 0 \end{cases} \quad (8.2.1)$$

来讨论。假定在方程组 (8.2.1) 解 x^* 的某邻域内， f_1 、 f_2 连续且有连续的一阶偏导数。现设 $(\xi_1^{(k)}, \xi_2^{(k)})$ 为上述邻域内的点，作线性函数

$$\begin{cases} l_1(\xi_1, \xi_2) = f_1(\xi_1^{(k)}, \xi_2^{(k)}) + \frac{\partial f_1}{\partial \xi_1} \Big|_{(\xi_1^{(k)}, \xi_2^{(k)})} (\xi_1 - \xi_1^{(k)}) \\ \quad + \frac{\partial f_1}{\partial \xi_2} \Big|_{(\xi_1^{(k)}, \xi_2^{(k)})} (\xi_2 - \xi_2^{(k)}), \\ l_2(\xi_1, \xi_2) = f_2(\xi_1^{(k)}, \xi_2^{(k)}) + \frac{\partial f_2}{\partial \xi_1} \Big|_{(\xi_1^{(k)}, \xi_2^{(k)})} (\xi_1 - \xi_1^{(k)}) \\ \quad + \frac{\partial f_2}{\partial \xi_2} \Big|_{(\xi_1^{(k)}, \xi_2^{(k)})} (\xi_2 - \xi_2^{(k)}), \end{cases} \quad (8.2.2)$$

在该邻域内近似替代函数 $f_1(\xi_1, \xi_2)$ 和 $f_2(\xi_1, \xi_2)$ 。于是，非线性方程组 (8.2.1) 即以线性方程组

$$\begin{cases} \frac{\partial f_1}{\partial \xi_1} \Big|_{x=x_k} \Delta \xi_1^{(k)} + \frac{\partial f_1}{\partial \xi_2} \Big|_{x=x_k} \Delta \xi_2^{(k)} = -f_1(x_k), \\ \frac{\partial f_2}{\partial \xi_1} \Big|_{x=x_k} \Delta \xi_1^{(k)} + \frac{\partial f_2}{\partial \xi_2} \Big|_{x=x_k} \Delta \xi_2^{(k)} = -f_2(x_k) \end{cases} \quad (8.2.3)$$

近似替代, 其中 $\Delta \xi_1^{(k)} = \xi_1 - \xi_1^{(k)}$, $\Delta \xi_2^{(k)} = \xi_2 - \xi_2^{(k)}$, $x_k = (\xi_1^{(k)}, \xi_2^{(k)})^T$. 如果在上述邻域中矩阵

$$\begin{pmatrix} \frac{\partial f_1}{\partial \xi_1} & \frac{\partial f_1}{\partial \xi_2} \\ \frac{\partial f_2}{\partial \xi_1} & \frac{\partial f_2}{\partial \xi_2} \end{pmatrix}_{x=x_k} \quad (8.2.4)$$

为非异, 则可解得

$$\begin{pmatrix} \Delta \xi_1^{(k)} \\ \Delta \xi_2^{(k)} \end{pmatrix} = \begin{pmatrix} \frac{\partial f_1}{\partial \xi_1} & \frac{\partial f_1}{\partial \xi_2} \\ \frac{\partial f_2}{\partial \xi_1} & \frac{\partial f_2}{\partial \xi_2} \end{pmatrix}_{x=x_k}^{-1} \begin{pmatrix} -f_1 \\ -f_2 \end{pmatrix}_{x=x_k}. \quad (8.2.5)$$

矩阵 (8.2.4) 常称为向量函数 $F(x)$ 的 Jacobi 矩阵, 记作 $F'(x)$. 进而, 记 $\xi_1^{(k+1)} = \xi_1^{(k)} + \Delta \xi_1^{(k)}$, $\xi_2^{(k+1)} = \xi_2^{(k)} + \Delta \xi_2^{(k)}$, 则 (8.2.5) 即可写作

$$\Delta x_k = -F'(x_k)^{-1} F(x_k)$$

或

$$x_{k+1} = x_k - F'(x_k)^{-1} F(x_k). \quad (8.2.6)$$

公式 (8.2.6) 称为求解非线性方程组 (8.2.1) 的牛顿迭代公式, 而线性方程组 (8.2.3) 则称为牛顿方程组.

我们对向量函数 $F(x)$ 自始至终作如下假设:

假设 1.

(1) 向量函数 $F(x)$ 在开的凸集 D^* 内连续可微, 即在 D 内函数 $f_i(\xi_1, \xi_2, \dots, \xi_n)$ ($i=1, \dots, n$) 连续且有连续的一阶偏导数 (请注意: 下文中集合 D 总是假定为开的凸集, 以后不再说明).

(2) D 内存在点 \mathbf{x}^* , 使 $\mathbf{F}(\mathbf{x}^*) = 0$ 且 $\mathbf{F}'(\mathbf{x}^*)$ 为非异, $\mathbf{F}'(\mathbf{x})$ 为 Jacobi 矩阵

$$\begin{pmatrix} \frac{\partial f_1}{\partial \xi_1} & \frac{\partial f_1}{\partial \xi_2} & \cdots & \frac{\partial f_1}{\partial \xi_n} \\ \frac{\partial f_2}{\partial \xi_1} & \frac{\partial f_2}{\partial \xi_2} & \cdots & \frac{\partial f_2}{\partial \xi_n} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial f_n}{\partial \xi_1} & \frac{\partial f_n}{\partial \xi_2} & \cdots & \frac{\partial f_n}{\partial \xi_n} \end{pmatrix} \quad (8.2.7)$$

综上所述, 牛顿法的基本思想是, 将非线性方程组逐次线性化, 从而形成迭代算法. 具体地说, 设 \mathbf{x}_k 为 \mathbf{x}^* 的第 k 次近似, 在 \mathbf{x}_k 的某邻域中, 如果线性函数

$$\mathbf{L}_k(\mathbf{x}) = \mathbf{F}(\mathbf{x}_k) + \mathbf{F}'(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) \quad (8.2.8)$$

是向量函数 $\mathbf{F}(\mathbf{x})$ 的一个较好的逼近, 则由线性方程组 $\mathbf{L}_k(\mathbf{x}) = 0$ 求得的解 \mathbf{x}_{k+1} 必然是 \mathbf{x}^* 更好的近似解.

牛顿法的计算步骤可总结成如下算法.

算法8.2.1 设 $\mathbf{F}: \mathbf{R}^n \rightarrow \mathbf{R}^n$, \mathbf{x}_0 为 \mathbf{x}^* 的初始近似, 它们都满足定理8.4.1或定理8.4.2的条件, 欲求满足一定精度的近似解.

- 1) 对 $k = 0, 1, \dots, m$ (m 为允许的最大迭代次数),
 - 1.1) 计算 $\mathbf{F}(\mathbf{x}_k)$, 如果 \mathbf{x}_k 满足精度要求, 则转 3).
 - 1.2) 计算 $\mathbf{F}'(\mathbf{x}_k)$.
 - 1.3) 求解线性方程组 $\mathbf{F}'(\mathbf{x}_k) \Delta \mathbf{x}_k = -\mathbf{F}(\mathbf{x}_k)$, 置 $\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta \mathbf{x}_k$.
- 2) NEXT k .
- 3) 置 $\mathbf{x}^* = \mathbf{x}_{k+1}$.
- 4) END.

牛顿法的特点是, 把向量函数 $\mathbf{F}(\mathbf{x})$ 看成一个整体, 通过

线性化，从而形成迭代算法，这种方法的优点是，便于收敛性分析，利于迭代格式的修正，但这种方法并未充分利用函数的具体结构。对于一个非线性向量函数 $F(\mathbf{x})$ ，其非线性程度在各个分量函数 $f_i(\mathbf{x})$ ($i = 1, \dots, n$) 上分布可能是不均衡的。也就是说，有的分量非线性程度高，有的则接近线性函数，甚至有的完全是线性函数。在这种情况下，把所有分量 $f_i(\mathbf{x})$ 采取完全相同的数值处理，将不利于方法整体计算效率的提高。于是人们设想：采取对 f_i ($i = 1, \dots, n$) 逐个线性化结合代入法以形成迭代算法的途径，这也就是 **Brown** 算法的基本思想。下面简要叙述 **Brown** 方法的迭代步骤。

设 \mathbf{x}_k 为 \mathbf{x}^* 的第 k 次近似。首先将 $f_1(\mathbf{x})$ 在 \mathbf{x}_k 近傍线性化，当 \mathbf{x}_k 充分接近 \mathbf{x}^* 时，即有

$$f_1(\mathbf{x}) \approx f_1(\mathbf{x}_k) + \sum_{j=1}^n \frac{\partial f_1(\mathbf{x}_k)}{\partial \xi_j} (\xi_j - \xi_j^{(k)}).$$

令上式为零，则可解出一个变量，设为 ξ_n 。则

$$\begin{aligned} \xi_n = \xi_n^{(k)} - \left(\frac{\partial f_1(\mathbf{x}_k)}{\partial \xi_n} \right)^{-1} \left\{ \sum_{j=1}^{n-1} \frac{\partial f_1(\mathbf{x}_k)}{\partial \xi_j} (\xi_j - \xi_j^{(k)}) \right. \\ \left. + f_1(\mathbf{x}_k) \right\}. \end{aligned} \quad (8.2.9)$$

上式右端是变量 $\xi_1, \xi_2, \dots, \xi_{n-1}$ 的线性函数。定义

$$\begin{aligned} L_n(\xi_1, \xi_2, \dots, \xi_{n-1}) = \xi_n^{(k)} - \left(\frac{\partial f_1(\mathbf{x}_k)}{\partial \xi_n} \right)^{-1} \\ \times \left\{ \sum_{j=1}^{n-1} \frac{\partial f_1(\mathbf{x}_k)}{\partial \xi_j} (\xi_j - \xi_j^{(k)}) + f_1(\mathbf{x}_k) \right\}. \end{aligned} \quad (8.2.10)$$

再考虑第二个分量函数 $f_2(\mathbf{x})$ ，把 $L_n(\xi_1, \xi_2, \dots, \xi_{n-1})$ 作为函数变量 ξ_n 的近似，把它代入 $f_2(\mathbf{x})$ ，这样 $f_2(\mathbf{x})$ 就变成一个 $n-1$ 个变量的多元函数，即

$$g_2(\mathbf{x}) \equiv f_2(\xi_1, \xi_2, \dots, \xi_{n-1}, L(\xi_1, \dots, \xi_{n-1})) \quad (8.2.11)$$

对 $g_2(\mathbf{x})$ 作类似于 $f_1(\mathbf{x})$ 的处理，即能解得

$$\xi_{n-1} = \xi_{n-1}^{(k)} - \left(\frac{\partial g_2(\mathbf{x}_k)}{\partial \xi_{n-1}} \right)^{-1} \left\{ \sum_{j=1}^{n-2} \frac{\partial g_2(\mathbf{x}_k)}{\partial \xi_j} (\xi_j - \xi_j^{(k)}) + g_2(\mathbf{x}_k) \right\}, \quad (8.2.12)$$

注意上式中 $\mathbf{x}_k = (\xi_1^{(k)}, \dots, \xi_{n-1}^{(k)})^T$. 上式右端为变量 $\xi_1, \xi_2, \dots, \xi_{n-1}$ 的线性函数, 记作 $L_{n-1}(\xi_1, \dots, \xi_{n-1})$.

按照这样的步骤, 依次处理 $F(\mathbf{x})$ 的各个分量函数. 最后考虑 $g_n(\mathbf{x})$, 此时 $g_n(\mathbf{x})$ 已简化为只含 ξ_1 的单变量函数. 通过线性化, 便能解得

$$\xi_1 = \xi_1^{(k)} - \left(\frac{dg_n(\mathbf{x}_k)}{d\xi_1} \right)^{-1} g_n(\mathbf{x}_k), \quad (8.2.13)$$

于是即求得了 ξ_1^* 的第 $k+1$ 个近似值. 由此出发利用线性函数 L_2, L_3, \dots, L_n 逐一回代, 最后将能求得新的迭代点 \mathbf{x}_{k+1} . 现举例说明以上迭代过程.

设非线性方程组为

$$\begin{cases} f_1(\mathbf{x}) \equiv \xi_1 \xi_3 - 2\xi_1 + 1 = 0, \\ f_2(\mathbf{x}) \equiv \xi_1 + 5\xi_2 + 1 = 0, \\ f_3(\mathbf{x}) \equiv \xi_1^2 + \xi_2^2 - \xi_3 - 2 = 0; \end{cases}$$

初始近似解为 $\mathbf{x}_0 = (-2, 0, 2)^T$.

根据前述步骤, 首先在点 \mathbf{x}_k 近傍用线性函数

$$f_1(\mathbf{x}_k) + (\xi_1 - \xi_1^{(k)}) (\xi_3^{(k)} - 2) + (\xi_3 - \xi_3^{(k)}) \xi_1^{(k)}$$

来近似替代 $f_1(\mathbf{x})$. 令上式为零, 即可解出变量

$$\xi_3 = \xi_3^{(k)} - (\xi_1^{(k)})^{-1} [(\xi_3^{(k)} - 2)\xi_1 + 1].$$

定义

$$L_3(\xi_1, \xi_2) = \xi_3^{(k)} - (\xi_1^{(k)})^{-1} [(\xi_3^{(k)} - 2)\xi_1 + 1].$$

下一步即可利用 $L_3(\xi_1, \xi_2)$ 消去 $f_2(\mathbf{x})$ 中的变量 ξ_3 (当然在本问题中, $f_2(\mathbf{x})$ 中并不含 ξ_3), 从而定义新函数

$$g_2(\xi_1, \xi_2) \equiv \xi_1 + 5\xi_2 + 1.$$

$g_2(\xi_1, \xi_2)$ 在 $(\xi_1^{(k)}, \xi_2^{(k)})$ 处的线性近似表达式为

$$g_2(\xi_1^{(k)}, \xi_2^{(k)}) + (\xi_1 - \xi_1^{(k)}) + 5(\xi_2 - \xi_2^{(k)}).$$

令上式为零, 即可解出变量

$$\xi_2 = -\frac{1}{5}(1 + \xi_1).$$

定义

$$L_2(\xi_1) \equiv -\frac{1}{5}(1 + \xi_1)$$

利用 $L_3(\xi_1, \xi_2)$ 和 $L_2(\xi_1)$ 消去 $f_3(x)$ 中的变量 ξ_2, ξ_3 , 从而定义新函数

$$g_3(\xi_1) \equiv \xi_1^2 + \frac{1}{25}(1 + \xi_1)^2 \\ - \{ \xi_2^{(k)} - (\xi_1^{(k)})^{-1} [(\xi_2^{(k)} - 2)\xi_1 + 1] \} - 2$$

$g_3(\xi_1)$ 在点 $\xi_1^{(k)}$ 近傍用线性函数

$$g_3(\xi_1^{(k)}) + \left(\frac{52}{25}\xi_1^{(k)} + \frac{\xi_2^{(k)} - 2}{\xi_1^{(k)}} + \frac{2}{25} \right) (\xi_1 - \xi_1^{(k)})$$

近似替代。令上式为零, 即可解得

$$\xi_1^{(k+1)} = \xi_1^{(k)} -$$

$$\frac{(\xi_1^{(k)})^2 + \frac{1}{25}(1 + \xi_1^{(k)})^2 - \{ \xi_2^{(k)} - (\xi_1^{(k)})^{-1} [(\xi_2^{(k)} - 2)\xi_1^{(k)} + 1] \} - 2}{\frac{52}{25}\xi_1^{(k)} + \frac{\xi_2^{(k)} - 2}{\xi_1^{(k)}} + \frac{2}{25}}$$

利用以上公式从 x_0 出发逐次迭代, 可得到如下结果:

n	ξ_1	ξ_2	ξ_3
0	-2.000000	0.000000	2.000000
1	-2.112745	0.222549	2.500000
2	-2.103978	0.220796	2.475893
3	-2.103937	0.220787	2.475299
4	-2.103937	0.220787	2.475299

综上所述, **Brown** 算法是对非线性方程组逐个线性化, 逐个消去变量以及向后回代相结合的迭代过程. 该算法在 §5 还将详细讨论.

在结束本节之前, 还有一个问题须作出回答. 前面我们将向量函数 $F(x)$ 的 **Jacobi** 矩阵用导数记号 $F'(x)$ 来表示, 那末, $F(x)$ 的 **Jacobi** 矩阵究竟如何体现出“导数”的意义呢?

回想单变量实值函数在点 x 处导数的概念. 如果存在一个实数 $a = f'(x)$, 使得

$$\lim_{h \rightarrow 0} \left(\frac{1}{h} \right) [f(x+h) - f(x) - ah] = 0$$

成立, 则 a 即为函数 $f(x)$ 在点 x 处的导数. 这个定义可以自然地推广到 n 维.

定义 8.2.1 设 $F: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$, x 为 D 的内点. 若存在 $A \in \mathbb{R}^{m \times n}$ 使对所有 $h \in \mathbb{R}^n$ 有

$$\lim_{h \rightarrow 0} \left(\frac{1}{\|h\|} \right) [\|F(x+h) - F(x) - Ah\|] = 0 \quad (8.2.14)$$

或

$$F(x+h) = F(x) + Ah + o(\|h\|). \quad (8.2.15)$$

成立, 则称 $F(x)$ 在点 x 处为 **Fréchet** 可微 (或 F -可微).

据范数等价定理知, (8.2.14) 中的极限与 \mathbb{R}^m 上所取的范数无关. 也就是说, 在某种范数意义下, 如果 $F(x)$ 在 x 处为 F -可微, 则对任何范数而言, $F(x)$ 在 x 处仍为 F -可微. 不难证明: 使 (8.2.14) 成立的矩阵 A 至多只有一个.

现设法用 F 的分量函数 f_1, \dots, f_m 的偏导数来表示矩阵 A . 设 $A = (a_{ij})$, 因为极限式 (8.2.14) 对任何趋于零的向量 $h \in \mathbb{R}^n$ 均成立, 故若取 $h = te_j$ ($t > 0$), 则

$$\lim_{t \rightarrow 0} \left(\frac{1}{t} \right) [f_i(x + te_j) - f_i(x) - ta_{ij}] = 0$$

$$(i = 1, \dots, m).$$

据假定, $f_i (i = 1, \dots, m)$ 在 \mathbf{x} 处的偏导数存在, 故上式意味着

$$a_{i,j} = \frac{\partial f_i}{\partial \xi_j} \quad (i = 1, \dots, m, j = 1, \dots, n).$$

由此可知, 在极限式 (8.2.14) 成立的前提下, $\mathbf{F}(\mathbf{x})$ 的 Jacobi 矩阵就是 (8.2.14) 式中的 \mathbf{A} 矩阵. 因此, 将 $\mathbf{F}(\mathbf{x})$ 的 Jacobi 矩阵记为 $\mathbf{F}'(\mathbf{x})$ 是极自然的. 当然, 反过来 Jacobi 矩阵的存在并不能保证 $\mathbf{F}(\mathbf{x})$ 为 F -可微. 例如, 当 $m=1$ 时, 即 $\mathbf{F}(\mathbf{x})$ 为 n 元函数时, 假定在点 \mathbf{x} 它的所有偏导数均存在, 那末它的 Jacobi 矩阵即为

$$\mathbf{F}'(\mathbf{x}) = \left(\frac{\partial F}{\partial \xi_1}, \frac{\partial F}{\partial \xi_2}, \dots, \frac{\partial F}{\partial \xi_n} \right), \quad (8.2.16)$$

但 $\mathbf{F}(\mathbf{x})$ 在 \mathbf{x} 点未必可微, 也就是说, $\mathbf{F}(\mathbf{x})$ 在点 \mathbf{x} 处未必 F -可微.

定义 8.2.2 设 $\mathbf{F}: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$, 如果 \mathbf{F} 在 D 的内点为 F -可微, 则把使 (8.2.14) 式成立的 $\mathbf{A} \in \mathbb{R}^{m \times n}$ 记为 $\mathbf{F}'(\mathbf{x})$, 即

$$\mathbf{F}'(\mathbf{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial \xi_1} & \frac{\partial f_1}{\partial \xi_2} & \dots & \frac{\partial f_1}{\partial \xi_n} \\ \frac{\partial f_2}{\partial \xi_1} & \frac{\partial f_2}{\partial \xi_2} & \dots & \frac{\partial f_2}{\partial \xi_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial f_m}{\partial \xi_1} & \frac{\partial f_m}{\partial \xi_2} & \dots & \frac{\partial f_m}{\partial \xi_n} \end{pmatrix}, \quad (8.2.17)$$

并称它为 $\mathbf{F}(\mathbf{x})$ 在 \mathbf{x} 处的 F -导数.

至此我们即可明了: 函数 $\mathbf{F}(\mathbf{x})$ 满足假设 1 (a) 实际上就是指 $\mathbf{F}(\mathbf{x})$ 在 D 内 F -可微. 事实上, 如果函数 $f_i(\xi_1, \dots, \xi_n)$ ($i = 1, \dots, m$) 在 D 内连续且有连续的一阶偏导数, 则对一切 i 成立

$$\lim_{\mathbf{h} \rightarrow 0} \frac{f_i(\mathbf{x} + \mathbf{h}) - f_i(\mathbf{x}) - \nabla f_i(\mathbf{x})^T \mathbf{h}}{\|\mathbf{h}\|} = 0, \quad (8.2.18)$$

因此

$$\lim_{h \rightarrow 0} \frac{F(x+h) - F(x) - F'(x)h}{\|h\|} = 0, \quad (8.2.19)$$

而 (8.2.19) 式等价于

$$\lim_{h \rightarrow 0} \left(\frac{1}{\|h\|} \right) \|F(x+h) - F(x) - F'(x)h\| = 0.$$

此即 $F(x)$ 在点 x 处 F -可微。

例8.2.1 用牛顿法求非线性方程组

$$\begin{cases} f_1(\xi_1, \xi_2) \equiv 4\xi_1^2 + \xi_2^2 + 2\xi_1\xi_2 - \xi_2 - 2 = 0, \\ f_2(\xi_1, \xi_2) \equiv 2\xi_1^2 + 3\xi_1\xi_2 + \xi_2^2 - 3 = 0, \end{cases}$$

的解, 初始近似取 $x_0 = (0.4, 0.9)^T$.

据算法8.2.1首先计算 $F(x_0)$:

$$f_1(\xi_1^{(0)}, \xi_2^{(0)}) = -0.73,$$

$$f_2(\xi_1^{(0)}, \xi_2^{(0)}) = -0.79,$$

然后计算 $F'(x_0)$, 为此求出

$$f'_{1\xi_1} = 8\xi_1 + 2\xi_2, \quad f'_{1\xi_2} = 2\xi_1 + 2\xi_2 - 1,$$

$$f'_{2\xi_1} = 4\xi_1 + 3\xi_2, \quad f'_{2\xi_2} = 3\xi_1 + 2\xi_2.$$

于是

$$F'(x_0) = \begin{pmatrix} 5.0 & 1.6 \\ 4.3 & 3.0 \end{pmatrix}.$$

因此线性方程组 $F'(x_0)\Delta x_0 = -F(x_0)$ 为

$$5.0\Delta\xi_1^{(0)} + 1.6\Delta\xi_2^{(0)} = 0.73,$$

$$4.3\Delta\xi_1^{(0)} + 3.0\Delta\xi_2^{(0)} = 0.79,$$

解得 $\Delta x = (0.114, 0.100)^T$, 故 $x_1 = x_0 + \Delta x_0 = (0.514, 1.000)^T$. 重复以上计算, 得 $F(x_1) = (0.084784, 0.070392)^T$, 而

$$F'(x_1) = \begin{pmatrix} 6.112 & 2.028 \\ 5.056 & 3.542 \end{pmatrix}.$$

再解线性方程组

$$6.112\Delta\xi_1^{(1)} + 2.028\Delta\xi_2^{(1)} = -0.084784,$$

$$5.056\Delta\xi_1^{(1)} + 3.542\Delta\xi_2^{(1)} = -0.070392,$$

得 $\Delta x_1 = (-0.013826, -0.000138)^T$, $x_2 = x_1 + \Delta x_1 = (0.500174, 0.999862)^T$. 进而 $F(x_2) = (0.000768, 0.000387)^T$.

在点 x_2 处的 Jacobi 矩阵为

$$\begin{pmatrix} 6.00116 & 2.000072 \\ 5.000282 & 3.500246 \end{pmatrix}.$$

解线性方程组

$$\begin{cases} 6.00116\Delta\xi_1^{(2)} + 2.000072\Delta\xi_2^{(2)} = -0.000768, \\ 5.000282\Delta\xi_1^{(2)} + 3.500246\Delta\xi_2^{(2)} = -0.000387, \end{cases}$$

得 $\Delta x_2 = (-0.000174, 0.000138)^T$, $x_3 = x_2 + \Delta x_2 = (0.500000, 1.000000)^T$. 最后, 我们发现 $F(x_3) = 0$. 也就是说经过三次迭代, 竟然得到了原方程组的精确解. 当然, 这是一个特例, 一般情形只能得到近似解. 但从这个例子可以看出, 牛顿法的收敛速度是比较快的. 在§4, 我们还将具体分析牛顿法的收敛性.

§3 迭代过程的收敛率

本节的主要任务是对迭代序列 $\{x_k\} \subset R^n$ 收敛速度的度量方法予以定义.

定义8.3.1 设 $\{x_k\} \subset R^n$ 是由某迭代算法所生成, 它收敛于 (8.1.1) 的解 x^* . 若对某向量范数 $\|\cdot\|$, 存在一与 k 无关的 $q \in (0, 1)$ 以及正整数 k_0 , 当 $k \geq k_0$ 时有

$$\|x_{k+1} - x^*\| \leq q \|x_k - x^*\| \quad (8.3.1)$$

成立，则称迭代序列 $\{x_k\}$ 至少是线性收敛的。

以上定义意味着，当 $k \geq k_0$ 时，每经一次迭代，误差至少以因子 q 缩小。一个实用的迭代算法至少应具有线性收敛速度。

一个新的、有价值的局部加速方法，它应该具有下述超线性收敛性。以后我们介绍的大多数算法，在一定条件下均具有超线性收敛的性质。

定义8.3.2 设 $\{x_k\} \subset R^n$ 由某迭代算法所生成，它收敛于(8.1.1)的解 x^* 。若对某种向量范数 $\|\cdot\|$ ，存在一收敛于零的实数列 $\{q_k\}$ ，使对任何 k 成立

$$\|x_{k+1} - x^*\| \leq q_k \|x_k - x^*\|, \quad (8.3.2)$$

则称 $\{x_k\}$ 超线性收敛，或者更确切地说是商超线性收敛，有时也说该算法是局部超线性收敛。

根据以上定义即可推知：只要当 $k \geq k_0$ 且 $x_k \neq x^*$ 时，则 $\{x_k\}$ 超线性收敛就意味着

$$\lim_{k \rightarrow \infty} \|x_{k+1} - x^*\| / \|x_k - x^*\| = 0. \quad (8.3.3)$$

以下定理是序列 $\{x_k\}$ 超线性收敛的一个必要条件，它对构造某些迭代算法的收敛准则是有用的。

定理8.3.1 设序列 $\{x_k\}$ 超线性收敛于 x^* ，当 $k \geq k_0$ 且 $x_k \neq x^*$ 时，则

$$\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = 1. \quad (8.3.4)$$

但一般而言，反之未必成立。

证 因为对任何 $k \geq k_0$ 均有

$$\left| \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} - \frac{\|x_k - x^*\|}{\|x_k - x^*\|} \right| \leq \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|},$$

即

$$\left| \frac{\|x_{k+1} - x_k\|}{\|x_k - x^*\|} - 1 \right| \leq \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|}.$$

据 (8.3.3) 即可推知 (8.3.4) 成立.

上述定理的逆命题未必成立, 我们给出一个反例.

例8.3.1 设序列 $\{x_k\}$ 定义为

$$\begin{cases} x_{2i-1} = \frac{1}{i!}, \\ x_{2i} = 2x_{2i-1}, \end{cases} \quad (i=1, 2, \dots)$$

又设 $x^* = 0$, 则上述序列收敛于 x^* , 但据序列定义, 显然它并不超线性收敛于 x^* . 然而

$$\frac{\|x_{k+1} - x_k\|}{\|x_k - x^*\|} = \begin{cases} 1 & (k=2i-1, i \geq 1), \\ 1 - \frac{1}{2(i+1)} & (k=2i, i \geq 1). \end{cases}$$

于是, 对该序列而言, (8.3.4) 是成立的.

定理8.3.1告诉我们, 如果已知 $\{x_k\}$ 超线性收敛于 x^* , 那末就可利用 $\|x_{k+1} - x_k\|$ 来估计误差 $\|x_k - x^*\|$. 据此, 即可构成如下迭代收敛准则: 事先指定正小数 ε_1 , 如果有

$$\|x_{k+1} - x_k\| \leq \varepsilon_1 \|x_k\| \quad (8.3.5)$$

或

$$\|x_{k+1} - x_k\| \leq \varepsilon_1 \|x_{k+1}\| \quad (8.3.6)$$

成立, 那末迭代终止. 有时候以上迭代收敛准则常和

$$\|F(x_k)\| < \varepsilon_2 \quad (8.3.7)$$

联用, 以防止出现只顾一头的偏向, 这里 ε_2 也为事先指定的正小数 (当然还可以采取其它措施). 这样做的目的是, 在实际计算时, $\{x_k\}$ 并不一定超线性收敛. 其次, 即便 $\{x_k\}$ 超线性收敛, $\|x_{k+1} - x_k\|$ 也未必为 k 的单调减函数.

§4 牛顿法的收敛性分析

牛顿法的主要优点可见于以下二个定理。

定理8.4.1 (牛顿法的局部收敛性)

设 $F: R^n \rightarrow R^n$ 满足假设 1, 则存在一闭球 $S(\bar{x}^*, \delta) \subset D$, 对任何 $x_0 \in S(\bar{x}^*, \delta)$, 牛顿迭代解

$$x_{k+1} = x_k - F'(x_k)^{-1}F(x_k)$$

对任何 k 均有定义且恒不越出该球。序列 $\{x_k\}$ 超线性收敛于 x^* 。

定理8.4.2 (牛顿法的平方收敛性)

设 $F: R^n \rightarrow R^n$ 满足假设 1, $F'(x)$ 在点 x^* 满足 Lipschitz 条件, 亦即存在正常数 K 使

$$\|F'(x) - F'(x^*)\| \leq K\|x - x^*\|, \quad \forall x \in D \quad (8.4.1)$$

成立, 则存在正常数 q 使

$$\|x_{k+1} - x^*\| \leq q\|x_k - x^*\|^2 \quad (k = 0, 1, \dots) \quad (8.4.2)$$

成立, 即 $\{x_k\}$ 至少是平方收敛的。

请注意: 和一维情形类似, 可以证明, 如果 D 充分小, 而 $F(x)$ 在 x^* 处二次可微*, 那末 (8.4.1) 即能满足。所谓 $F(x)$ 在 x 处二次可微, 这里是指 $F(x)$ 的所有分量函数 $f_i(x)$ ($i = 1, \dots, n$) 在 x 处二次连续可微。

*) 仿照定义 8.2.1, 对二阶导数可作如下定义:

设 $F: D \subset R^n \rightarrow R^n$ 在 D 内 Fréchet 可微, x 为 D 的内点, 若存在双线性算子 $A: R^n \times R^n \rightarrow R^n$ 使

$$\lim_{h \rightarrow 0} \left(\frac{1}{\|h\|} \right) \|F'(x+h) - F'(x) - Ah\| = 0$$

或

$$F'(x+h) = F'(x) + Ah + o(\|h\|)$$

对所有 $h \in R^n$ 成立, 则 $F(x)$ 在点 x 处二次可微并称算子 A 为 F 在点 x 处的二阶导数, 记 $A = F''(x)$ 。

今后我们称 $F(x)$ 在 D 上满足 (8.4.1) 为假设2.

为了证明以上二个定理, 我们首先证明下面几个引理.

引理8.4.1 设矩阵 $A(x) \in R^{n \times n}$ 在点 $x^* \in D$ 处连续且 $A(x^*)$ 为非异, 则存在闭球 $\overline{S(x^*, \delta)} \subset D$ 及正数 γ , 使对一切 $x \in \overline{S(x^*, \delta)}$, $A(x)$ 亦为非异且

$$\|A(x)^{-1}\| \leq \gamma.$$

证 令 $\alpha = \|A(x^*)^{-1}\|$, 取 $\beta < \frac{1}{\alpha}$, 对取定的 β 选取这样的 $\delta > 0$, 以保证 $\overline{S(x^*, \delta)} \subset D$ 且当 $x \in \overline{S(x^*, \delta)}$ 时有

$$\|A(x^*) - A(x)\| \leq \beta$$

成立. 据 $A(x)$ 在 x^* 处的连续性, 这是能办到的. 考虑到当 $x \in \overline{S(x^*, \delta)}$ 时有

$$\|I - A(x^*)^{-1}A(x)\| \leq \|A(x^*)^{-1}\| \|A(x^*) - A(x)\| \leq \alpha\beta < 1, \quad (8.4.3)$$

又

$$A(x^*)^{-1}A(x) = I - (I - A(x^*)^{-1}A(x)),$$

因此, 当 $x \in \overline{S(x^*, \delta)}$ 时, $A(x^*)^{-1}A(x)$ 必为非异, 于是 $A(x)$ 亦必为非异. 又当 $x \in \overline{S(x^*, \delta)}$ 时应有

$$\begin{aligned} \|A(x)^{-1}\| &= \|[I - (I - A(x^*)^{-1}A(x))]^{-1}A(x^*)^{-1}\| \\ &\leq \|A(x^*)^{-1}\| \|[I - (I - A(x^*)^{-1}A(x))]^{-1}\| \\ &\leq \|A(x^*)^{-1}\| \frac{1}{1 - \|I - A(x^*)^{-1}A(x)\|} \\ &\leq \frac{\alpha}{1 - \alpha\beta} = \gamma. \end{aligned}$$

有时以上引理常以下面形式出现: 设 $A_k \in R^{n \times n}$ 而 $A \in R^{n \times n}$ 非异且

$$\|A_k - A\| \leq \beta,$$

若 $\|A^{-1}\| \leq \alpha$ 且 $\alpha\beta < 1$, 则 A_k 亦为非异且

$$\|A_k^{-1}\| \leq \frac{\alpha}{1 - \alpha\beta}. \quad (8.4.4)$$

文献上常称这个引理为 **Banach** 引理, 证明过程则完全和引理 8.4.1 相同. **Banach** 引理的意义就是: 若 A_k 能接近一个可逆阵, 那末 A_k 本身也是可逆的.

引理 8.4.2 设 $\varphi: [a, b] \subset \mathbb{R}^1 \rightarrow \mathbb{R}^n$ 连续, 则有

$$\left\| \int_a^b \varphi(t) dt \right\| \leq \int_a^b \|\varphi(t)\| dt. \quad (8.4.5)$$

证: 由于 $\|\varphi(t)\|$ 是 t 的连续函数, 故 $\|\varphi(t)\|$ 黎曼可积. 据向量函数的积分定义知, 对任意给定的 $\varepsilon > 0$, 存在分划 $a = t_0 < \dots < t_m = b$ 使

$$\left\| \int_a^b \varphi(t) dt - \sum_{i=1}^m \varphi(t_i) (t_i - t_{i-1}) \right\| \leq \varepsilon. \quad (8.4.6)$$

另一方面, 据一元函数积分定义有

$$\left| \int_a^b \|\varphi(t)\| dt - \sum_{i=1}^m \|\varphi(t_i)\| (t_i - t_{i-1}) \right| \leq \varepsilon. \quad (8.4.7)$$

由 (8.4.6) 和 (8.4.7) 即可推知

$$\begin{aligned} \left\| \int_a^b \varphi(t) dt \right\| &\leq \left\| \sum_{i=1}^m \varphi(t_i) (t_i - t_{i-1}) \right\| + \varepsilon \\ &\leq \sum_{i=1}^m \|\varphi(t_i)\| (t_i - t_{i-1}) + \varepsilon \\ &\leq \int_a^b \|\varphi(t)\| dt + 2\varepsilon, \end{aligned}$$

因为 ε 为任意小的正数, 故 (8.4.5) 成立.

中值定理在一元函数或多元函数微积分中占有十分重要的地位. 对于向量函数的微积分情况也是这样. 不过, 它的中值定理的形式略有变化.

引理 8.4.3 设 $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$ 满足假设 1(1), 则对任意 $z, y \in D$ 有不等式

$$\|F(y) - F(x)\| \leq \max_{0 \leq t \leq 1} \|F'[x + t(y-x)]\| \|y-x\| \quad (8.4.8)$$

成立。

证 由于 D 为开的凸集, 故对 $0 \leq t \leq 1$, $x + t(y-x) \in D$, 因此 $F'(x + t(y-x))$ 有定义. 对任何 $y, x \in D$ 以下等式显然成立.

$$\begin{aligned} f_i(y) - f_i(x) &= \int_0^1 \frac{d}{dt} f_i[\xi_1 + t(\eta_1 - \xi_1), \dots, \xi_n + t(\eta_n - \xi_n)] dt \\ &= \int_0^1 \sum_{j=1}^n \frac{\partial}{\partial \xi_j} f_i[x + t(y-x)] (\eta_j - \xi_j) dt \\ &= \int_0^1 f'_i[x + t(y-x)] (y-x) dt, \end{aligned}$$

于是得到 $F(x)$ 积分形式的中值定理:

$$F(y) - F(x) = \int_0^1 F'[x + t(y-x)] (y-x) dt, \quad (8.4.9)$$

两端取范数并利用不等式 (8.4.5), 即可推得 (8.4.8).

引理 8.4.4 设 $F: R^n \rightarrow R^n$ 满足假设 1(1), 则对任何 $x, y, z \in D$ 有以下不等式

$$\begin{aligned} \|F(y) - F(z) - F'(x)(y-z)\| &\leq \max_{0 \leq t \leq 1} \|F'[x + t(y-z)] \\ &\quad - F'(x)\| \|y-z\| \end{aligned} \quad (8.4.10)$$

成立. 进而, 如果 $F'(x)$ 在 D 中 Lipschitz 连续, 即对任何 $x, y \in D$, 存在正常数 κ , 使

$$\|F'(x) - F'(y)\| \leq \kappa \|x - y\| \quad (8.4.11)$$

成立, 则

$$\|F(y) - F(z) - F'(x)(y-z)\| \leq \kappa \sigma(y, z) \|y-z\|, \quad (8.4.12)$$

其中 $\sigma(y, z) = \max\{\|y-x\|, \|z-x\|\}$. 特别地, 当 $z=x$ 时有不等式

$$\|F(y) - F(x) - F'(x)(y-x)\| \leq \frac{1}{2} \kappa \|y-x\|^2 \quad (8.4.13)$$

成立, 这里 κ 为 Lipschitz 常数.

证 利用 (8.4.9) 和引理 8.4.2 可推得

$$\begin{aligned}\|F(y) - F(x) - F'(x)(y - x)\| &= \left\| \int_0^1 \{F'[x + t(y - x)] - F'(x)\} (y - x) dt \right\| \\ &\leq \int_0^1 \|F'[x + t(y - x)] - F'(x)\| \|y - x\| dt \\ &\leq \max_{0 \leq t \leq 1} \|F'[x + t(y - x)] - F'(x)\| \|y - x\|.\end{aligned}$$

如果 $F'(x)$ 在 D 中 Lipschitz 连续, 那末可进一步推得

$$\begin{aligned}\|F(y) - F(x) - F'(x)(y - x)\| &\leq \kappa \int_0^1 \|x + t(y - x) - x\| \|y - x\| dt \\ &\leq \kappa \max_{0 \leq t \leq 1} \|x + t(y - x) - x\| \|y - x\| \\ &= \kappa \sigma(y, x) \|y - x\|.\end{aligned}$$

而

$$\begin{aligned}\|F(y) - F(x) - F'(x)(y - x)\| &\leq \kappa \|y - x\|^2 \int_0^1 t dt \\ &= \frac{1}{2} \kappa \|y - x\|^2.\end{aligned}$$

至此即能证明早已提出的定理 8.4.1 和定理 8.4.2.

定理 8.4.1 (牛顿法的局部收敛性) 的证明:

据定理 8.4.1 假定, $F'(x)$ 在点 $x^* \in D$ 处连续且 $F'(x^*)$ 为非异, 故利用引理 8.4.1 知: 存在闭球 $S(x^*, \delta') \subset D$ 及正数 γ , 对一切 $x \in S(x^*, \delta')$, $F'(x)$ 为非异且 $\|F'(x)^{-1}\| \leq \gamma$.

另一方面, 据 $F'(x)$ 在 x^* 处的连续性知, 对任给定的 $\varepsilon > 0$ 总存在 $\delta > 0$ ($\delta \leq \delta'$), 对一切 $x \in S(x^*, \delta)$ 有

$$\|F'(x) - F'(x^*)\| \leq \varepsilon$$

成立. 又因 $F(x)$ 于 x^* 处 F -可微, 我们可将 δ 选取充分小, 使 $2\varepsilon\gamma = q < 1$ 且不等式

$$\|F(x) - F(x^*) - F'(x^*)(x - x^*)\| \leq \varepsilon \|x - x^*\|$$

$$\forall x \in \overline{S(x^*, \delta)}$$

成立。

现设 $x \in S(x^*, \delta)$ 并记 $G(x) = x - F'(x)^{-1}F(x)$ ，则有

$$\begin{aligned} \|G(x) - x^*\| &= \|x - x^* - F'(x)^{-1}[F(x) - F(x^*)]\| \\ &\leq \| -F'(x)^{-1} \| [\|F(x) - F(x^*) \\ &\quad - F'(x^*)(x - x^*)\| \\ &\quad + \|F'(x) - F'(x^*)\| \|x - x^*\|] \\ &\leq 2\varepsilon \gamma \|x - x^*\| = q \|x - x^*\|. \end{aligned} \quad (8.4.14)$$

现若取 x_0 使 $\|x_0 - x^*\| \leq \delta$ ，则据 (8.4.14) 且记 $x_1 = G(x_0)$ ，有

$$\|x_1 - x^*\| \leq q \|x_0 - x^*\| \leq q\delta < \delta.$$

由此可见 $x_1 \in \overline{S(x^*, \delta)}$ 。利用归纳法即可证

$$\|x_k - x^*\| \leq q \|x_{k-1} - x^*\| < \delta \quad (8.4.15)$$

对所有的 $k \geq 1$ 均成立，这里 $x_k = G(x_{k-1})$ 。逐次利用这些不等式即可得到

$$\|x_k - x^*\| \leq q^k \|x_0 - x^*\| \quad (k = 1, 2, \dots), \quad (8.4.16)$$

又因 $0 < q < 1$ ，故

$$\lim_{k \rightarrow \infty} \|x_k - x^*\| = 0,$$

亦即

$$\lim_{k \rightarrow \infty} x_k = x^*.$$

这样，我们就证明了牛顿法的局部收敛性。

再令

$$\begin{aligned} \gamma \left[\frac{\|F(x_k) - F(x^*) - F'(x^*)(x_k - x^*)\|}{\|x_k - x^*\|} \right. \\ \left. + \|F'(x_k) - F'(x^*)\| \right] = q_k, \end{aligned}$$

则据 (8.4.14) 即有

$$\|x_{k+1} - x^*\| \leq q_k \|x_k - x^*\|, \quad (k = 0, 1, \dots)$$

据 $F(x)$ 在 x^* 的 F -可微性以及 $F'(x)$ 在 x^* 处的连续性即可推知

$$\lim_{k \rightarrow \infty} q_k = 0,$$

故牛顿迭代序列 $\{x_k\}$ 是局部超线性收敛的。

定理 8.4.2 的证明:

由不等式

$$\|x_{k+1} - x^*\| \leq \|F'(x_k)^{-1}\| \|F(x_k) - F(x^*) - F'(x_k)(x_k - x^*)\|,$$

并用引理 8.4.4 即可推得

$$\|x_{k+1} - x^*\| \leq \frac{3}{2} \gamma \kappa \|x_k - x^*\|^2 \quad (k = 0, 1, \dots),$$

由此可见 $\{x_k\}$ 至少是平方收敛。

至此, 我们有必要对定理 8.4.1 的证明予以回顾。牛顿迭代法实际上可看作形如

$$x = G(x) \quad (8.4.17)$$

的非线性方程组的简单迭代法, 即

$$x_{k+1} = G(x_k) \quad (k = 0, 1, \dots). \quad (8.4.18)$$

定理 8.4.1 主要论证了函数 $G(x)$ 满足以下三个条件:

(1) 存在一闭球 $S(x^*, \delta)$, 函数 $G(x)$ 在该闭球上有定义;

(2) $G(S(x^*, \delta)) \subset S(x^*, \delta)$;

(3) 存在 $q \in (0, 1)$ 使

$$\|G(x) - G(x^*)\| \leq q \|x - x^*\|, \quad \forall x \in \overline{S(x^*, \delta)}. \quad (8.4.19)$$

第三个条件通常称为**压缩条件**。在定理 8.4.1 中这个条件实际上是

$$G'(x^*) = 0 \quad (8.4.20)$$

或

$$\|G'(\mathbf{x}^*)\| = 0,$$

这可从 (8.4.14) 看出。一般地, 若 $G(\mathbf{x})$ 在 \mathbf{x}^* 处连续可微且

$$\rho(G'(\mathbf{x}^*)) < 1, \quad (8.4.21)$$

即能保证存在球 $S(\mathbf{x}^*, \delta)$, $G(\mathbf{x})$ 在 $S(\mathbf{x}^*, \delta)$ 上满足压缩条件。事实上, 据范数理论知: 对任意给定的 $\varepsilon > 0$, 存在某种范数 $\|\cdot\|_k$ 使

$$\|G'(\mathbf{x}^*)\|_k \leq \rho(G'(\mathbf{x}^*)) + \varepsilon, \quad (8.4.22)$$

而对 $\varepsilon > 0$ 又存在 $S(\mathbf{x}^*, \delta)$, 当 $\mathbf{x} \in S(\mathbf{x}^*, \delta)$ 时有

$$\begin{aligned} \|\varphi(\mathbf{x}) - G(\mathbf{x}^*)\|_k &\leq \int_0^1 \|G'[\mathbf{x}^* + t(\mathbf{x} - \mathbf{x}^*)]\|_k \|\mathbf{x} - \mathbf{x}^*\|_k dt \\ &\leq (\|G'(\mathbf{x}^*)\|_k + \varepsilon) \|\mathbf{x} - \mathbf{x}^*\|_k \end{aligned} \quad (8.4.23)$$

成立, 从而

$$\|G(\mathbf{x}) - G(\mathbf{x}^*)\|_k \leq (\rho(\varphi'(\mathbf{x}^*)) + 2\varepsilon) \|\mathbf{x} - \mathbf{x}^*\|_k.$$

由于 ε 的任意性, 故总有 $q = (\rho(\varphi'(\mathbf{x}^*)) + 2\varepsilon) < 1$ 成立, 即 $G(\mathbf{x})$ 在 $S(\mathbf{x}^*, \delta)$ 上满足压缩条件。

从以上几个定理可以看出: 牛顿法有二大优点 (Brown 方法亦类似)。其一, 在一般条件下牛顿法是超线性收敛的; 如果再加上假设 2 这个条件, 那末牛顿法至少是平方收敛的。粗略地说, 后者意味着牛顿法每迭代一次, 有效数位基本上增加了一倍, 这就是说, 如果由准确到一位的近似解出发, 那只要经四次迭代就可得到准确到十六位的近似解。其二, 牛顿法在一般条件下总存在一个吸引域 S 。换言之, 只要初始近似解落在 S 中, 那么由此出发, 生成的迭代序列 $\{\mathbf{x}_k\}$ 总不会越出域 S , 而且收敛于 \mathbf{x}^* 。

另外也应指出: 牛顿法是自动校正的, 这就是说 \mathbf{x}_{k+1} 仅仅依赖于 F 和 \mathbf{x}_k , 因而前面迭代时产生的舍入误差不会一步

步传下去。牛顿法这一优点我们将会看到并不能为拟牛顿法所共有。

正如我们在引言中早就指出的，确定一个初始近似解 x_0 有时是很困难的，这正是牛顿法的一大缺点（Brown方法也一样）。因为对某一具体问题，吸引区域 S 很可能极小，因此初始近似解必须选得很好，然而这常常是很难办到的。克服这一缺点的措施之一，可以在牛顿公式中引入收敛因子 $\lambda_k > 0$ ，以扩大收敛区域。这时迭代公式具有如下形式：

$$x_{k+1} = x_k - \lambda_k F'(x_k)^{-1} F(x_k) \quad (k = 0, 1, \dots) \quad (8.4.24)$$

而 λ_k 如此选择，以保证 $\{\|F(x_k)\|\}$ 为严格单调下降序列。选取的办法通常使 $\lambda_k > 0$ 为极值问题

$$\min_{\lambda > 0} \|F[x_k - \lambda F'(x_k)^{-1} F(x_k)]\| \quad (8.4.25)$$

的解。设 $\varphi(\lambda) \equiv \|F[x_k - \lambda F'(x_k)^{-1} F(x_k)]\|$ 为区间 $[a, b]$ 上的严格单峰函数^{*}，则极小点 λ_k 可采用0.618优选法寻找。0.618搜索方法的步骤如下：

在 λ_k 的预估区间 $[a, b]$ 中，计算两个试验点 $a_1 = a + 0.381966(b - a)$ ， $b_1 = a + 0.618034(b - a)$ 以及其上的函数值 φ_1 和 φ_2 。如果 $\varphi_1 < \varphi_2$ ，则取 $[a, b_1]$ 为新的估值区间，否则取 $[a_1, b]$ 。然后重复上述过程。由于试验点取法特殊，下一次要计算的两点之一可由前次算出试验点近似替代，而另外一点又近似等于区间端点之和减去已知的试验点。所以，除第一次需要计算两个点上的函数值外，以后只要计算一个新点上的函数

* 一个函数 $\varphi: [a, b] \subset \mathbb{R}^1 \rightarrow \mathbb{R}^1$ 在 $[a, b]$ 上是严格单峰的，如果存在 $\lambda^* \in [a, b]$ 为 $\varphi(\lambda)$ 在 $[a, b]$ 上的整体极小点，并且对任何 $\lambda', \lambda'' \in [a, b]$ ，当 $\lambda' < \lambda''$ 时，若 $\lambda'' \leq \lambda^*$ ，则 $\varphi(\lambda') > \varphi(\lambda'')$ ；若 $\lambda^* \leq \lambda'$ ，则 $\varphi(\lambda'') > \varphi(\lambda')$ 。

值。当区间长度比原始区间长度小到0.01倍时，搜索结束并取其中点作为 λ_k 。 λ_0 的预估区间可取 $[0, 1.618034\lambda_{-1}]$ ，其中 λ_{-1} 事先选定，可取0.1或1。 λ_k 的预估区间可取 $[0, 1.618034\lambda_{k-1}]$ ，这里 λ_{k-1} 为前次迭代的收敛因子。

牛顿法另一个缺点 (Brown 方法也类似) 是：对每个 k 都必须计算矩阵 $F'(x_k)$ ，向量 $F(x_k)$ 以及 $O(n^3)$ 次算术运算求解线性方程组。这就是说，每一步都要计算 n^2 个偏导函数值。对于大多数 $F(x)$ 而言，完成它将花费很大的工作量。因此，如果对某一问题的 **Jacobi** 矩阵相对而言计算量不太大，如 **Jacobi** 矩阵为大型稀疏矩阵，那末牛顿法是值得采用的；如果对某一问题的 **Jacobi** 矩阵的计算量很大，那末可选用各种修正的牛顿法或其它方法。例如，为了减少牛顿法中矩阵求逆的次数，一个较为普遍采用的技巧是：在一定的迭代次数内，使 **Jacobi** 矩阵固定，这对于 **Jacobi** 矩阵变化很慢的情形是十分有效的。具体说，定义 $x_{k+1} = y_{k,k}$ ， $y_{k,0} = x_k$ 而

$$y_{k,j} = y_{k,j-1} - F'(x_k)^{-1} F(x_{k,j-1}) \quad (j = 1, 2, \dots, k_0),$$

(8.4.26)

可以证明由 (8.4.26) 产生的迭代序列 $\{x_k\}$ 收敛阶*至少为 $k_0 + 1$ (习题8.7)，所以它也是产生高阶迭代过程的一个有效且简单的途径。

最后还必须指出：如果 $F'(x^*)$ 奇异，那末牛顿法的收敛性质就要被破坏。近年来，对于奇异点的牛顿法已有许多研究成果。

●) 设序列 $\{x_k\}$ 收敛于 x^* ，则 $\{x_k\}$ 的收敛阶为

$$p = \sup\{q \in \mathbb{R}^+ \mid \lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^q} \in [0, \infty) \text{ 而 } q \geq 1\}.$$

§5 离散的牛顿法

由前节分析我们知道，牛顿法的主要缺点是对每个 k 都必须计算矩阵 $F'(x_k)$ ，亦即对每个 k 都必须计算 n^2 个导数。因此，人们很自然地会想到：矩阵 $F'(x_k)$ 中 (i, j) 元用一个差商来近似替代。换言之，将 Jacobi 矩阵离散化。通常最简单的办法是， $\frac{\partial f_i}{\partial \xi_j}$ 以一阶差商

$$\frac{f_i(x_k + h_k e_j) - f_i(x_k)}{h_k}$$

替代之。其中 e_j 表示第 j 个单位向量， h_k 常取作 $O(\|F(x_k)\|)$ 或 $O(|f_1(x_k)|)$ 。在此基础上可构成一类离散型牛顿算法 N_k 。

算法 8.5.1 (N_k 算法)

1) 对 $k = 0, 1, \dots, m$ 。

1.1) 对充分小的步长 h_k ，计算矩阵 J_k ，它的第 j 列为

$$J_k e_j = \frac{F(x_k + h_k e_j) - F(x_k)}{h_k} \quad (j = 1, \dots, n), \quad (8.5.1)$$

1.2) 用近似 Jacobi 矩阵 J_k (假定它为非异) 执行 k_0 次“牛顿迭代”。亦即定义 $x_{k+1} = y_{k,k_0}$ ， $y_{k,0} = x_k$ 而

$$y_{k,j} = y_{k,j-1} - J_k^{-1} F(y_{k,j-1}) \quad (j = 1, \dots, k_0), \quad (8.5.2)$$

1.3) 如果 x_{k+1} 满足精度要求，则转 3)。

1.4) 计算 $F(x_{k+1})$ 以及 $F(x_{k+1} + h_{k+1} e_j)$ ($j = 1, \dots, n$)。

2) NEXT k 。

3) 置 $x^* = x_{k+1}$ 。

4) END。

与习题8.7类似，在一定条件下，可以证明 N_k 算法的收敛阶至少为 $k_0 + 1$ 。顺便指出， N_k 算法每迭代一次需计算 $n + k_0$ 次函数值，（这里是指向量函数 $F(x)$ 的函数值，其每次计算应包括 n 次纯量函数值）。

很自然 Brown 算法同样也能离散化，离散 Brown 算法的主要步骤如下：

第一步：设 x_k 为 (8.1.1) 解 x^* 的近似值，则对函数 f_1 在 x_k 的近傍用线性函数

$$f_1(x_k) + f_{1\xi_1;h_k}(x_k)(\xi_1 - \xi_1^{(k)}) + \cdots + f_{1\xi_n;h_k}(x_k)(\xi_n - \xi_n^{(k)}) \quad (8.5.3)$$

近似替代。其中 $x = (\xi_1, \dots, \xi_n)^T$ ， $x_k = (\xi_1^{(k)}, \dots, \xi_n^{(k)})^T$ 而

$$f_{1\xi_j;h_k}(x_k) = \frac{f_1(x_k + h_j^{(k)} e_j) - f_1(x_k)}{h_j^{(k)}}. \quad (8.5.4)$$

（以下为书写简便计，略去了 h 的上下标，请读者注意）。如果 x_k 充分接近 x^* ，那末可令 (8.5.3) 式为零，并由此解出一个变量。譬如说，如果 $f_{1\xi_n;h_k}(x_k)$ 按模最大，则就可解出 ξ_n 来，即

$$\xi_n = \xi_n^{(k)} - \sum_{j=1}^{n-1} (f_{1\xi_j;h}^{(k)} / f_{1\xi_n;h}^{(k)}) (\xi_j - \xi_j^{(k)}) - f_1^{(k)} / f_{1\xi_n;h}^{(k)}, \quad (8.5.5)$$

其中 $f_{1\xi_j;h}^{(k)} \equiv f_{1\xi_j;h_k}(x_k)$ ， $f_1^{(k)} \equiv f_1(x_k)$ 。可以证明， $F(x)$ 在假设 1 的条件下，至少有一个非零偏导数 $\frac{\partial f_1(x_k)}{\partial \xi_j}$ 。因此，相应的近似偏导数 $f_{1\xi_j;h}^{(k)}$ 亦为非零。这样 (8.5.5) 总是有意义的。选取绝对值为最大的近似偏导数相除，其作用类似于用高斯消去法求解线性方程组时，采用部分主元素方案，以保证方法的数值稳定性。由 (8.5.5) 式可以看出， ξ_n 乃是 $n-1$ 个变量 $\xi_1, \xi_2, \dots, \xi_{n-1}$ 的线性函数。为书写清楚起见，记 (8.5.5) 的左端为 $L_n(\xi_1, \dots, \xi_{n-1})$ 且定义 $L_n^{(k)} \equiv L_n(\xi_1^{(k)}, \dots, \xi_{n-1}^{(k)})$ 。

第二步：对函数 f_1 ，定义新函数

$$g_2(\xi_1, \dots, \xi_{n-1}) \equiv f_2(\xi_1, \dots, \xi_{n-1}, L_n(\xi_1, \dots, \xi_{n-1})) \quad (8.5.6)$$

且记 $g_2^{(k)} \equiv f_2(\xi_1^{(k)}, \dots, \xi_{n-1}^{(k)}, L_n^{(k)})$. 注意这时 g_2 是变量 ξ_1, \dots, ξ_{n-1} 的函数. 同样地, 在点 $(\xi_1^{(k)}, \dots, \xi_{n-1}^{(k)})$ 近傍以线性函数近似替代 g_2 , 并由此解出一个变量. 设 ξ_{n-1} 对应的近似偏导数 $g_{2\xi_{n-1}h}$ 按绝对值为最大, 于是解得 ξ_{n-1} 为:

$$\xi_{n-1} = \xi_{n-1}^{(k)} - \sum_{j=1}^{n-2} (g_{2\xi_j h}^{(k)} / g_{2\xi_{n-1} h}^{(k)}) (\xi_j - \xi_j^{(k)}) - g_2^{(k)} / g_{2\xi_{n-1} h}^{(k)}, \quad (8.5.7)$$

这里近似偏导数 $g_{2\xi_j h}^{(k)}$ 为

$$g_{2\xi_j h}^{(k)} \equiv \frac{g_2(\xi_1^{(k)}, \dots, \xi_j^{(k)} + h, \xi_{j+1}^{(k)}, \dots, \xi_{n-1}^{(k)}) - g_2^{(k)}}{h}. \quad (8.5.8)$$

由 (8.5.7) 知, ξ_{n-1} 为其余 $n-2$ 个变量的线性函数. 现将此线性函数记作 $L_{n-1}(\xi_1, \dots, \xi_{n-2})$, 同样定义 $L_{n-1}^{(k)} \equiv L_{n-1}(\xi_1^{(k)}, \dots, \xi_{n-2}^{(k)})$. 请注意: 此时应把 L_n 中的 ξ_{n-1} 换为 L_{n-1} , 即此时 L_n 的变量为 $\xi_1, \dots, \xi_{n-2}, L_{n-1}(\xi_1, \dots, \xi_{n-2})$.

继续进行这种“逐步代入”的过程, 至第 $i+1$ 步, 得到线性表达式

$$L_{n-i} \equiv \xi_{n-i}^{(k)} - \sum_{j=1}^{n-i-1} (g_{i+1, \xi_j h}^{(k)} / g_{i+1, \xi_{n-i} h}^{(k)}) (\xi_j - \xi_j^{(k)}) - g_{i+1}^{(k)} / g_{i+1, \xi_{n-i} h}^{(k)}, \quad (8.5.9)$$

其中

$$g_{i+1}^{(k)} \equiv f_{i+1}(\xi_1^{(k)}, \dots, \xi_{n-i}^{(k)}, L_{n-i+1}^{(k)}, \dots, L_n^{(k)}), \quad (8.5.10)$$

$$g_{i+1, \xi_j h}^{(k)} \equiv \frac{g_{i+1}(\xi_1^{(k)}, \dots, \xi_j^{(k)} + h, \xi_{j+1}^{(k)}, \dots, \xi_{n-i}^{(k)}) - g_{i+1}^{(k)}}{h}$$

$$(j = 1, \dots, n-i). \quad (8.5.11)$$

由此可见, 在完成第 $i+1$ 步时, 一个主要的工作就是要计算 g_{i+1} 在点 $(\xi_1^{(k)}, \dots, \xi_{n-i}^{(k)})^T \equiv \mathbf{x}_{n-i}^{(k)}$ 以及点 $\mathbf{x}_{n-i}^{(k)} + h\mathbf{e}_j$ ($j = 1, \dots, n-i$)

处的值。请注意：虽然计算 g_{i+1} 就是计算 f_{i+1} ，但这时 f_{i+1} 自变量的后 i 个分量为 $L_n, L_{n-1}, \dots, L_{n-i+1}$ 。而 $L_n, L_{n-1}, \dots, L_{n-i+1}$ 的自变量又分别为：

$$L_n: (\xi_1, \dots, \xi_{n-i}, L_{n-i+1}, \dots, L_{n-1});$$

$$L_{n-1}: (\xi_1, \dots, \xi_{n-i}, L_{n-i+1}, \dots, L_{n-2});$$

$$\vdots \quad \quad \quad \vdots$$

$$L_{n-i+2}: (\xi_1, \dots, \xi_{n-i}, L_{n-i+1})$$

$$L_{n-i+1}: (\xi_1, \dots, \xi_{n-i}).$$

由此可见， $L_n, L_{n-1}, \dots, L_{n-i+1}$ 在点 $x_{n-i}^{(k)}$ 和 $x_{n-i}^{(k)} + h e_j$ ($j = 1, \dots, n-i$) 的值，应自 L_{n-i+1} 出发通过逐步回代求得。回代方程取

$$L_{s+1} = \xi_{s+1}^{(k)} - \sum_{j=1}^s (g_{n-s, \xi_j, h} / g_{n-s, \xi_{s+1}, h}) (L_j - \xi_j^{(k)}) \\ - g_{n-s} / g_{n-s, \xi_{s+1}, h} \quad (s = n-i, \dots, n-1), \quad (8.5.12)$$

(注意：在计算点 $x_{n-i}^{(k)}$ 处的值时，实际上 j 只要从 $n-i+1$ 开始)。其次，在确定了 L_{n-i} 之后，在 $L_n, L_{n-1}, \dots, L_{n-i+1}$ 中凡出现 ξ_{n-i} 之处均用 L_{n-i} 替代。

第 n 步：这时有

$$g_n \equiv f_n(\xi_1, L_2, \dots, L_n), \quad (8.5.13)$$

其中 L_i ($i = 2, 3, \dots, n$) 应具有如下形式：

$$L_i = \xi_i^{(k)} - \sum_{j=1}^{i-1} (g_{n-i+1, \xi_j, h} / g_{n-i+1, \xi_i, h}) (L_j - \xi_j^{(k)}) \\ - g_{n-i+1} / g_{n-i+1, \xi_i, h}. \quad (8.5.14)$$

它们当然都是变量 ξ_1 的函数（为完整起见，定义 $L_1 \equiv \xi_1$ ， $g_1 \equiv f_1$ ），因此 g_n 亦为单变量 ξ_1 的函数。在点 $\xi_1^{(k)}$ 近傍，以线性函数来近似替代 g_n ，并由此解出

$$\xi_1 = \xi_1^{(k)} - g_n^{(k)} / g_{n, \xi_1, h}^{(k)}, \quad (8.5.15)$$

于是求得 $\xi_i^{(k+1)}$, 即 ξ_i^* 的第 $k+1$ 个近似值. 现记 $L_i \equiv \xi_i^{(k+1)}$ 并利用三角方程组 (8.5.14) 逐步回代, 即可解得 $\xi_1^{(k+1)}, \dots, \xi_n^{(k+1)}$.

按以上步骤反复迭代, 直到满足迭代收敛条件 (8.3.5) 或 (8.3.6) 和 (8.3.7) 为止.

至此, 我们将上述方法简要总结成如下算法.

算法 8.5.2 (Brown 方法)

1) 对 $k = 0, 1, \dots, m$ (m 为允许的最大迭代次数).

1.1) 对 $i = 1, 2, \dots, n$ (n 为未知数的个数),

1.1.1) 形成第 i 个函数 g_i 的各个离散偏导数 $g_{i,j,k}$,

1.1.2) 寻找其中绝对值为最大的偏导数, 即

$$|g_{i,j_0,k}| = \max_j \{|g_{i,j,k}|\},$$

1.1.3) 形成方程组 (8.5.14) 中第 i 个方程 (关于 L_i) 的系数.

1.2) NEXT i .

1.3) 回代以求得下一个近似解 x_{k+1} .

1.4) 如满足收敛条件则转 3).

2) NEXT k .

3) END.

牛顿法每迭代一次要计算 n^2 个偏导数以及计算 n 个分量函数值. 那末, Brown 方法每迭代一次又需要计算多少次函数值呢?

据以上算法可知, $i=1$ 时, 需计算 f_1 的 $n+1$ 个值; $i=2$ 时, 需计算 n 个 f_1 的值, 依次类推. 由此可知, Brown 方法每次迭代需

$$\sum_{i=1}^{n+1} i = \frac{n^2 + 3n}{2}$$

次函数值计算。从这里显然可以看出“逐步代入”法的优点。当然应强调指出：这里函数值计算的节省仅指原始方程组 (8.1.1) 中 f_i 而言。事实上，Brown 方法还包括好些其它函数的计算，即线性函数 L_k 的计算。但这些计算并没有计入 $\frac{1}{2}n^2 + \frac{3}{2}n$ 之中。因此，只有当函数 f_i 的计算较繁复时，Brown 方法相对于 N_1 算法而言，计算量才有真正的节省。

因为 Brown 方法每次只对一个方程 $f_i = 0$ 运算，而在处理下一个方程 $f_{i+1} = 0$ 时，又即刻利用已得到的信息。因此，对于方程组 (8.1.1) 就存在一个最优次序的问题。也就是说，方程应预先作这样的最优排列：线性方程排在最前，以后方程按接近线性的程度依次排列。能够采取这种最优次序的措施，是 Brown 方法的一个优点。

下面我们不加证明地引入 Brown 方法的局部收敛定理。

定理 8.5.1 设 $F: R^n \rightarrow R^n$ 满足假设 1，则必存在一闭球 $\overline{S(x^*, \delta)}$ 和正数 ε ，当 $x_0 \in \overline{S(x^*, \delta)}$ 且 $\{h_j^{(k)}\}$ 以 ε 为界时，由 Brown 方法生成的迭代序列 $\{x_k\}$ 收敛于 x^* 。如果 F 还满足假设 2 且 $h_j^{(k)} = O(1f_1(x_k))$ ，则 $\{x_k\}$ 至少是平方收敛的。

第 八 章 习 题

8.1 利用牛顿法和 Brown 方法求下列线性方程组

$$\text{I) } \begin{cases} \xi_1^2 + \xi_2^2 = 4, \\ \xi_1^2 - \xi_2^2 = 1, \end{cases} \quad \text{II) } \begin{cases} \xi_1^2 + \xi_2^2 + \xi_3^2 = 1, \\ 2\xi_1^2 + \xi_2^2 - 4\xi_3^2 = 0, \\ 3\xi_1^2 - 4\xi_2^2 + \xi_3^2 = 0, \end{cases}$$

准确到四位的近似解。初始近似解分别取 $x_0 = (1.6, 1.2)^T$ 和 $x_0 = (0.5, 0.5, 0.5)^T$ 。

8.2 设

$$f(x) = \begin{cases} 0, & x = 0, \\ [\xi_2(\xi_1^2 + \xi_2^2)^{\frac{3}{2}}]/[(\xi_1^2 + \xi_2^2)^2 + \xi_2^3], & x \neq 0, \end{cases}$$

证明 $f(x)$ 在 $x=0$ 处偏导数存在, 但不存在 F -导数.

8.3 设 $f(x) = \xi_1^2 + \xi_1 \xi_2$, 试按定义求 $f(x)$ 的一阶和二阶 Fréchet 导数.

8.4 试证, 使 (8.2.14) 成立的矩阵 A 至多只有一个.

8.5 设 $F: R^n \rightarrow R^n$ 在开凸集 D 上 F -可导, 试证 $F'(x)$ 在 D 内连续的充要条件是所有的偏导数 $\frac{\partial f_i}{\partial \xi_j}$ ($i, j = 1, 2, \dots, n$) 在 D 内连续.

8.6 设 $\{x_k\} \subset R^n$ 是由迭代算法 M 所生成, 它收敛于 (8.1.1) 的解 x^* . 若对某向量范数 $\|\cdot\|$ 成立等式

$$\lim_{k \rightarrow \infty} \|x_k - x^*\|^{1/k} = 0,$$

则称 $\{x_k\}$ R -超线性收敛于 x^* . 试证, 如果 $\{x_k\}$ 超线性收敛于 x^* , 则 $\{x_k\}$ 必 R -超线性收敛, 反之未必成立.

8.7 设 $A(x) \in R^{n \times n}$ 在点 $x^* \in D$ 处连续且 $A(x^*)$ 为非异, 则存在一闭球 $S(x^*, \delta) \subset D$, 在 $S(x^*, \delta)$ 上 $A(x)^{-1}$ 在 x^* 处连续.

8.8 设 $F: D \subset R^n \rightarrow R^n$ 在 x^* 处 F -可微, 定义在 D 上的矩阵 $A(x) \in R^{n \times n}$ 在 x^* 处连续且 $A(x^*)$ 非异, 则存在闭球 $S(x^*, \delta) \subset D$, 使向量函数

$$G(x) = x - [A(x)]^{-1}F(x)$$

在其上有定义且 $G(x)$ 在 x^* 处 F -可微, F -导数为

$$G'(x^*) = I - (A(x^*))^{-1}F'(x^*).$$

8.9 设 $F: R^n \rightarrow R^n$ 满足假设 1 和假设 2, 试证由 (8.4.26) 确定的迭代序列 $\{x_k\}$ 收敛阶至少为 $k_0 + 1$.

8.10 设非线性方程组为

$$f_1(\xi_1, \xi_2) \equiv \xi_1^2 + 10\xi_1\xi_2 + 4\xi_2^2 + 0.7401006 = 0,$$

$$f_2(\xi_1, \xi_2) \equiv \xi_1^2 - 3\xi_1\xi_2 + 2\xi_2^2 - 1.0201228 = 0,$$

试用带收敛因子的牛顿法求方程组的近似解. 初始近似解取 $x_0 = (0.1, -0.1)^T$, $\varepsilon_1 = \varepsilon_2 = 10^{-6}$ 或自定. 要求最后打印方程组近似解, 函数 f_1 ,

f_2 的值以及迭代次数。

8.11 设 $F: R^n \rightarrow R^n$ 满足假设 1, 试证存在一闭球 $\overline{S(x^*, \delta)}$ 以及正数 β 使

$$\|F(x)\| \geq \beta \|x - x^*\|, \quad \forall x \in \overline{S(x^*, \delta)}$$

成立。

8.12 编制一个离散 Brown 算法的程序, 算法中

$h_j^{(k)} (j=1, \dots, n)$ 按如下规定选取:

$$h_j^{(k)} = \max\{\alpha_j^{(k)}, 5 \times 10^{-\beta+2}\},$$

$$\alpha_j^{(k)} = \min\{\max(|f_1^{(k)}|, |g_1^{(k)}|, \dots, |g_n^{(k)}|), 0.001 \times |\xi_j^{(k)}|\},$$

其中 β 为机器的有效数位。

8.13 利用上题程序求以下非线性方程组的近似解:

$$i) \begin{cases} f_1(\xi_1, \xi_2) \equiv 2\xi_1^2 - \xi_2^2 - 1 = 0, \\ f_2(\xi_1, \xi_2) \equiv \xi_1^2 - \xi_2 - 4 = 0, \end{cases} \quad x_0 = (1.2, 1.7)^T$$

$$ii) \begin{cases} f_1(\xi_1, \xi_2) \equiv \cos(0.4\xi_2 + \xi_1^2) + \xi_1^2 + \xi_2^2 - 1.6 = 0, \\ f_2(\xi_1, \xi_2) \equiv 1.5\xi_1^2 - \frac{\xi_2^2}{0.36} - 1 = 0. \end{cases} \quad x_0 = (1.04, 0.47)^T,$$

参 考 文 献

- [1] Dennis, J.E & Moré, J.J. (1977). "Quasi-Newton Methods, Motivation and Theory" SIAM Review. vol. 19, No. 1 January.
- [2] Ortega, J. M. & W. C. Rheinboldt (1970), *Iterative Solution of Nonlinear Equations in Several variables*, Academic Press, New York.
- [3] Byrne, G. D. & Hall, C. A. (1972), *Numerical Solution of Systems of Nonlinear Algebraic Equations* Academic press New York.
- [4] 王德人 非线性方程组解法与最优化方法。
- [5] Березин, И. С. И Жидков, Н. П. (1966), *Методы Вычислений*.
- [6] Wait, R. (1979), *The Numerical Solution of Algebraic Equations*.

第九章 非线性方程组的拟 牛顿迭代解法

§1 拟牛顿方法的基本思想

在第八章§4我们曾指出：牛顿法(离散牛顿法亦类似)的缺点之一是，每迭代一次需计算 $n^2 + n$ 个纯量函数以及 $O(n^3)$ 次算术运算。现在随着实际问题维数的提高，迭代次数的增多，如何减少每次迭代所需的工作量就显得十分必要了。

1959、1965年，Davidon和Broyden针对无约束极小化问题以及非线性方程组(先后)提出了一类新的方法。1963年，Fletcher和Powell对Davidan的工作做了重要的修改和澄清，从而将这方面的研究工作引向了深入。这类新方法的基本思想是，用矩阵 B_k 近似替代 $F'(x_k)$ ，而 B_{k+1} 可在 B_k 的基础上用一个低秩矩阵来校正，所需的运算量为 $O(n)^2$ 次算术运算。显然，这样的方案计算量大大降低了。

当然计算量降低也是要付出代价的，偿还的代价就是将平方收敛速度降为超线性收敛速度。

这类新方法名称繁多，有拟牛顿法、变尺度方法、修正方法以及校正方法等等。下面我们针对非线性方程组来引入这类方法。

设 $F: R^n \rightarrow R^n$ 满足假设 1, $x_k \in D$ 而 $s_k \neq 0$, $x_k + s_k = x_{k+1} \in D$, 则据引理8.4.4知：对给定的 $\varepsilon > 0$, 总存在 $\delta > 0$, 当 $\|x_{k+1} - x_k\| < \delta$ 时使

$$\|F(x_k) - F(x_{k+1}) - F'(x_{k+1})(x_k - x_{k+1})\| \leq \varepsilon \|x_k - x_{k+1}\| \quad (9.1.1)$$

成立。因此有

$$F(x_k) \approx F(x_{k+1}) + F'(x_{k+1})(x_k - x_{k+1}), \quad (9.1.2)$$

其近似程度随 $\|x_k - x_{k+1}\|$ 的减小而增大。如果要求 B_{k+1} 来近似替代 $F'(x_{k+1})$ 的话, 那末很自然地要求 B_{k+1} 满足方程

$$F(x_k) = F(x_{k+1}) + B_{k+1}(x_k - x_{k+1}) \quad (9.1.3)$$

或

$$B_{k+1} s_k = F(x_{k+1}) - F(x_k). \quad (9.1.4)$$

如果 $F(x)$ 为单变量函数, 那末 (9.1.4) 即为

$$B_{k+1} = \frac{F(x_{k+1}) - F(x_k)}{x_{k+1} - x_k}.$$

上式右端表示点 x_{k+1} 处 $F(x)$ 的差商, 所以 (9.1.4) 常称为 **广义差商条件**。因为 B_{k+1} 可以看作 $F'(x_{k+1})$ 的近似矩阵, 所以 (9.1.4) 又称为拟牛顿方程, 它是发展拟牛顿法的一个中心。

前面已经指出: 除要求 B_{k+1} 满足拟牛顿方程之外, 我们还希望 B_{k+1} 能表示为 $B_k + \Delta B_k$, 这里 ΔB_k 为一低秩矩阵。这样就得到了一类迭代算法:

$$\begin{cases} x_{k+1} = x_k - B_k^{-1} F(x_k), \\ B_{k+1}(x_{k+1} - x_k) = F(x_{k+1}) - F(x_k), \\ B_{k+1} = B_k + \Delta B_k, \text{ rank } \Delta B_k = m \geq 1. \end{cases} \quad (k = 0, 1, \dots), \quad (9.1.5)$$

如果对所有的 k , $B_k^{-1} = H_k$ 存在, 那末 (9.1.5) 的过程是有意义的, 利用 Sherman-Morrison-Woodbury 公式 (习题 9.2) 可得到 ΔH_k 的一个显表达式。具体地说, 因为 $\text{rank}(\Delta B_k) = m$, 而任何秩为 m 的 n 阶矩阵都可分解为

$$\begin{aligned} \Delta B_k &= U_k V_k^T, \quad U_k, V_k \in R^{n \times m}, \\ \text{rank } U_k &= \text{rank } V_k = m. \end{aligned}$$

如果 $(I + V_k^T B_k^{-1} U_k)$ 为非异, 则有

$$(B_k + U_k V_k^T)^{-1} = B_k^{-1} - B_k^{-1} U_k (I + V_k^T B_k^{-1} U_k)^{-1} V_k^T B_k^{-1},$$

于是

$$\Delta H_k = Y_k W_k^T, \quad (9.1.6)$$

其中

$$Y_k = -B_k^{-1} U_k (I + V_k^T B_k^{-1} U_k)^{-1}, \quad (9.1.7)$$

$$W_k = (B_k^{-1})^T V_k. \quad (9.1.8)$$

显然 $Y_k, W_k \in R^{n \times m}$ 且它们的秩均为 m , 故 ΔH_k 仍为一个秩 m 矩阵. 这样 (9.1.5) 又可改写为

$$\begin{aligned} x_{k+1} &= x_k - H_k F(x_k), \\ H_{k+1} [F(x_{k+1}) - F(x_k)] &= x_{k+1} - x_k \end{aligned} \quad (k=0, 1, \dots) \quad (9.1.9)$$

$$H_{k+1} = H_k + \Delta H_k, \quad \text{rank } \Delta H_k = m.$$

以下很自然地要讨论 $\Delta B_k, \Delta H_k$ 的具体计算, 下节要介绍的 Broyden 方法, 就是一种秩1校正公式.

§2 Broyden 方法

为书写方便起见, 以下暂将矩阵和向量的下标删去, 并采用如下符号:

$$\bar{x} = x_{k+1}, \quad x = x_k, \quad s = \bar{x} - x, \quad y = F(\bar{x}) - F(x),$$

$$\bar{B} = B_{k+1}, \quad \bar{H} = H_{k+1}, \quad B = B_k, \quad H = H_k,$$

$$\Delta B = \bar{B} - B, \quad \Delta H = \bar{H} - H$$

$$u = y - Bs, \quad v = -Hu = s - Hy.$$

现在考虑校正矩阵 ΔB 为秩1阵的情形. 因为任何秩1阵均可表为二向量的外积, 故设

$$\Delta B = bc^T \quad b, c \in R^n \quad (9.2.1)$$

或

$$\bar{B} = B + bc^T \quad (9.2.2)$$

代入拟牛顿方程即得

$$(B + bc^T)s = y,$$

故得

$$u = \langle c, s \rangle b.$$

于是秩1校正公式应在下面范围中选取

$$\Delta B = \sigma uc^T,$$

而

$$\sigma = \begin{cases} 1/\langle c, s \rangle, & s \neq 0; \\ 0, & s = 0, \end{cases}$$

或

$$\bar{B} = B + \frac{(y - Bs)c^T}{\langle c, s \rangle} \quad (s \neq 0). \quad (9.2.3)$$

因为我们把 \bar{B} 看作 $F'(\bar{x})$ 的一个近似, 所以确定 \bar{B} 的条件应是 s 方向上的导数即拟牛顿方程, 而在与 s 正交的方向上我们未曾提出过什么要求。所以, 在与 s 正交的所有方向 z 上要求成立

$$\bar{B}z = Bz, \quad \forall z \in Z \equiv \{z | z \neq 0, \langle z, s \rangle = 0\} \quad (9.2.4)$$

这种想法是合乎逻辑的。将(9.2.3)和(9.2.4)结合起来, 即可得到唯一的校正公式

$$\Delta B = \sigma us^T, \\ \sigma = \begin{cases} 1/\langle s, s \rangle, & s \neq 0 \\ 0, & s = 0 \end{cases}$$

或

$$\bar{B} = B + \frac{(y - Bs)s^T}{\langle s, s \rangle}. \quad (s \neq 0). \quad (9.2.5)$$

这样, 我们就得到了 Broyden 秩 1 迭代算法。

算法9.2.1 设 $F: R^n \rightarrow R^n$, x_0 为 x^* 的初始近似解, B_0 为初始迭代矩阵(通常 B_0 就取作 $F'(x_0)$), 它们都符合定理9.3.1的假定.

1) 计算 $F(x_0)$.

2) 对 $k = 0, 1, \dots, m$,

2.1) 求解线性方程组,

$$B_k \Delta x_k = -F(x_k)$$

并置 $x_{k+1} = x_k + \Delta x_k$ ($\Delta x_k \equiv s_k$),

2.2) 计算 $F(x_{k+1})$, 如果 x_{k+1} 满足精度要求, 那末转4),

2.3) 计算 $y_k = F(x_{k+1}) - F(x_k)$,

2.4) 计算 B_{k+1}

$$B_{k+1} = B_k + \frac{(y_k - B_k s_k) s_k^T}{\langle s_k, s_k \rangle}.$$

3) NEXT k .

4) 置 $x^* = x_{k+1}$.

5) END.

Broyden 校正公式具有一个十分重要的性质, 即在所有满足拟牛顿方程的矩阵集合中, \bar{B} 是在 **Frobenius** 范数意义下最接近 B 的矩阵. 实际上, 这一性质等价于确定 \bar{B} 的条件(9.2.4).

定理9.2.1 设给定 $B \in R^{n \times n}$, $y \in R^n$ 以及非零向量 $s \in R^n$, \bar{B} 为由(9.2.5)确定的矩阵, 则 \bar{B} 为以下问题的唯一解.

$$\min\{\|\bar{B} - B\| \mid \bar{B}s = y\}. \quad (9.2.6)$$

证 利用(9.2.5)并考虑到 $y = \bar{B}s$ 及习题9.1有

$$\|\bar{B} - B\|_F = \left\| (\bar{B} - B) \frac{ss^T}{\langle s, s \rangle} \right\|_F$$

$$\leq \| \bar{B} - B \|_F = \left\| \frac{ss^T}{\langle s, s \rangle} - I \right\|_F$$

$$= \| \bar{B} - B \|_F,$$

所以 \bar{B} 确为问题(9.2.6)的解. 至于唯一性可以这样说明: 设 $f(A) = \|A - B\|_F$, 不难验证函数 $f(A)$ 在集合

$$N = \{A \mid As = y\} \quad (9.2.7)$$

上是一个严格凸函数*. 另一方面, 集合 N 显然是 $R^{n \times n}$ 中的一个凸集. 因此, \bar{B} 为凸函数 $f(A)$ 在凸集 N 上的唯一解(习题9.4).

对校正公式(9.2.5)利用 Sherman-Morrison 公式可得

$$\bar{H} = H - \frac{Hbs^T H}{1 + s^T Hb},$$

其中 $b = \sigma u$. 因为 $\sigma Hy = \sigma s + Hb$, 故 $1 + s^T Hb = \sigma s^T Hy$. 将它们代入上式即可得

$$\bar{H} = H - \frac{Hus^T H}{s^T Hy} = H + \frac{(s - Hy)s^T H}{s^T Hy}. \quad (9.2.9)$$

由此可见, 如果 $B^{-1} = H$ 存在, 那末, \bar{B} 为非异的充要条件为 $s^T Hy \neq 0$.

算法9.2.2 设 $F: R^n \rightarrow R^n$, x_0 为 x^* 的初始近似解, H_0 为初始迭代矩阵(通常 H_0 取作 $F'(x_0)^{-1}$, 它们都符合定理9.3.1的假设.

- 1) 计算 $F(x_0)$ 和 $F'(x_0)$.
- 2) 对 $k = 0, 1, \dots, m$,

*函数 $f: R^n \rightarrow R$ 在凸集 $D \subset R^n$ 上是凸的, 意指对任意 $x', x'' \in D$ 且 $x' \neq x''$ 以及任意的 $\lambda \in (0, 1)$ 成立

$$f[\lambda x' + (1 - \lambda)x''] \leq \lambda f(x') + (1 - \lambda)f(x'') \quad (9.2.8)$$

如果“ \leq ”改为“ $<$ ”, 那末函数 f 在 D 上是严格凸的.

- 2.1) 计算 $\mathbf{x}_{k+1} = \mathbf{x}_k - H_k \mathbf{F}(\mathbf{x}_k)$ ($-H_k \mathbf{F}(\mathbf{x}_k) = \mathbf{s}_k$),
- 2.2) 计算 $\mathbf{F}(\mathbf{x}_{k+1})$, 如果 \mathbf{x}_{k+1} 满足精度要求, 那末转4),
- 2.3) 计算 $\mathbf{y}_k = \mathbf{F}(\mathbf{x}_{k+1}) - \mathbf{F}(\mathbf{x}_k)$,
- 2.4) 计算 H_{k+1}

$$H_{k+1} = H_k + \frac{(\mathbf{s}_k - H_k \mathbf{y}_k) \mathbf{s}_k^T H_k}{\mathbf{s}_k^T H_k \mathbf{y}_k}.$$

3) NEXT k .

4) 置 $\mathbf{x}^* = \mathbf{x}_{k+1}$.

5) END.

比较算法 9.2.1 和算法 9.2.2, 当给定 \mathbf{x}_0 和 B_0 后, 算法 9.2.1 每迭代一次, 要计算 n 个纯量函数, 然后求解线性方程组, 因此每迭代一次还需 $O(n^3)$ 次算术运算. 算法 9.2.2 每迭一次亦需计算 n 个纯量函数. 另外, 形成校正公式 (9.2.9) 需 $O(n^2)$ 次算术运算. 由此可见, 采用算法 9.2.2 较好. 可是, 利用校正公式 (9.2.5) 也有可能迭代一次只需 $O(n^2)$ 次算术运算. 1972 年, Gill 和 Murray 提出了一种计算方案: 设已知 $B_k = Q_k R_k$ (这里 Q_k 为正交阵, R_k 为上三角阵), 则形成 B_{k+1} 的 QR 分解也只要 $O(n^2)$ 次算术运算. 当然, 如果 $B_k = Q_k R_k$, 那末, 求解线性方程组 $B_k \mathbf{s}_k = -\mathbf{F}(\mathbf{x}_k)$ 就只需 $O(n^2)$ 次算术运算. 这个计算方案要比算法 9.2.2 好, 这是因为校正公式 (9.2.5) 中毋需执行矩阵与向量的乘法, 因为 $B_k \mathbf{s}_k$ 即为 $-\mathbf{F}(\mathbf{x}_k)$. 另外, 理论分析也表明, 校正公式 (9.2.5) 更为稳定.

例 9.2.1 利用 Broyden 秩 1 方法求方程组

$$f_1(\xi_1, \xi_2) \equiv \xi_1^2 - \xi_2 - 1 = 0,$$

$$f_2(\xi_1, \xi_2) \equiv (\xi_1 - 2)^2 + (\xi_2 - 0.5)^2 - 1 = 0$$

的近似解. 初始近似解取 $\mathbf{x}_0 = (0, 0)^T$.

解 采用算法9.2.2计算. 首先计算 $F(x_0) = (-1, 3, 25)^T$.
因为

$$F'(x_0) = \begin{pmatrix} 2\xi_1 & -1 \\ 2\xi_1 - 4 & 2\xi_2 - 1 \end{pmatrix},$$

故

$$F'(x_0) = \begin{pmatrix} 0 & -1 \\ -4 & -1 \end{pmatrix}, \quad r_0 = F'(x_0)^{-1} = \begin{pmatrix} 0.25 & -0.25 \\ -1 & 0 \end{pmatrix}.$$

据牛顿法求得, $x_1 = (1.0625, -1)^T$; $F(x_1) = (1.12890625, 2.12890625)^T$; 按语句2.3)和2.4)算得

$$s_1 = \begin{pmatrix} 1.0652 \\ -1 \end{pmatrix}, \quad y_1 = \begin{pmatrix} 2.12890625 \\ -1.12109375 \end{pmatrix}$$

以及

$$H_1 = \begin{pmatrix} 0.3557441\dots & -0.2721932\dots \\ -0.5224991\dots & -0.1002162\dots \end{pmatrix}.$$

按语句2.1)算得

$$x_2 = \begin{pmatrix} 1.0625 \\ -1 \end{pmatrix} - \begin{pmatrix} 0.3557441\dots & -0.2721932\dots \\ -0.5224991\dots & -0.1002162\dots \end{pmatrix} \begin{pmatrix} 1.12890625 \\ 2.12890625 \end{pmatrix} = \begin{pmatrix} 1.240372062\dots \\ -0.1967974670\dots \end{pmatrix}.$$

经过11次迭代可得到具有12位有效数字的近似解:

$$x_{11} = (1.54634288332, 1.39117631279)^T.$$

如果用牛顿法,那末只要经七次迭代即能得到这一结果.因次,就一般而言, **Broyden** 方法的收敛速度要较牛顿方法收敛速度慢.当然,对两个联立方程的情形, **Broyden** 方法的优点显得不突出,甚至不如用牛顿法方便.只有当 n 较大时,在计算上才能显出它的优点.得到这一好处的代价就是收敛速度较牛顿法要稍慢一点.下一节我们就来分析 **Broyden** 方法的收敛性.

§3 Broyden方法的收敛性分析

首先应该强调指出：以下采用的证明技巧也适用于其它拟牛顿方法的收敛性分析。现假定 F 满足假设1和假设2，则当初始近似解 x_0 和初始矩阵 B_0 分别充分接近 x^* 和 $F'(x^*)$ 时，算法 9.2.1 生成的迭代序列 $\{x_k\}$ 超线性收敛于方程组的解 x^* 。为了证明这一结论，首先证明以下几个引理。

引理9.3.1 设 $F: R^n \rightarrow R^n$ 满足假设1和假设2, $x \in D$ 而 \bar{x} 由算法9.2.1生成且设 $\bar{x} \in D$, 则

$$\|\bar{B} - F'(x^*)\|_p \leq \|B - F'(x^*)\|_p + \kappa \sigma(x, \bar{x}), \quad (p=2, F) \quad (9.3.1)$$

其中

$$\sigma(x, \bar{x}) = \max\{\|\bar{x} - x^*\|, \|x - x^*\|\}. \quad (9.3.2)$$

证 由(9.2.5)有

$$\bar{B} - F'(x^*) = (B - F'(x^*)) \left(I - \frac{s s^T}{\langle s, s \rangle} \right) + \frac{(y - F'(x^*)s) s^T}{\langle s, s \rangle}, \quad (9.3.3)$$

上式两端取范数，据习题9.1~9.3即有

$$\|\bar{B} - F'(x^*)\|_p \leq \|B - F'(x^*)\|_p + \frac{\|y - F'(x^*)s\|}{\|s\|}.$$

据引理 8.4.4 有

$$\|y - F'(x^*)s\| \leq \kappa \sigma(x, \bar{x}) \|s\|$$

成立。至此(9.3.1)得证。

引理9.3.2 设 $F: R^n \rightarrow R^n$ 满足假设1和假设2, 又设 $\{B_k\}$ 为一非异矩阵序列。如果对某 $x_0 \in D$ 迭代序列

$$\mathbf{x}_{k+1} = \mathbf{x}_k - B_k^{-1} \mathbf{F}(\mathbf{x}_k) \quad (k = 0, 1, \dots)$$

恒在 D 中且 $\mathbf{x}_k \neq \mathbf{x}^*$ ($\forall k \leq 0$)，又设该迭代序列收敛于 \mathbf{x}^* ，那末当且仅当

$$\lim_{k \rightarrow \infty} \frac{\| [B_k - \mathbf{F}'(\mathbf{x}^*)](\mathbf{x}_{k+1} - \mathbf{x}_k) \|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|} = 0 \quad (9.3.4)$$

时， $\{\mathbf{x}_k\}$ 超线性收敛于 \mathbf{x}^* 。

证 若 (9.3.4) 成立，则因

$$\begin{aligned} [B_k - \mathbf{F}'(\mathbf{x}^*)](\mathbf{x}_{k+1} - \mathbf{x}_k) &= -\mathbf{F}(\mathbf{x}_k) - \mathbf{F}'(\mathbf{x}^*)(\mathbf{x}_{k+1} - \mathbf{x}_k) \\ &= [\mathbf{F}(\mathbf{x}_{k+1}) - \mathbf{F}(\mathbf{x}_k) - \mathbf{F}'(\mathbf{x}^*)(\mathbf{x}_{k+1} - \mathbf{x}_k)] \\ &\quad - \mathbf{F}(\mathbf{x}_{k+1}), \end{aligned} \quad (9.3.5)$$

于是

$$\begin{aligned} \lim_{k \rightarrow \infty} \frac{\|\mathbf{F}(\mathbf{x}_{k+1})\|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|} &\leq K \lim_{k \rightarrow \infty} \sigma(\mathbf{x}_k, \mathbf{x}_{k+1}) \\ &\quad + \lim_{k \rightarrow \infty} \frac{\|[B_k - \mathbf{F}'(\mathbf{x}^*)](\mathbf{x} - \mathbf{x}_k)\|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|}, \end{aligned}$$

故

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{F}(\mathbf{x}_{k+1})\|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|} = 0. \quad (9.3.6)$$

又因 $\mathbf{F}'(\mathbf{x}^*)$ 为非异，故存在一正数 β ，当 k 充分大以后有

$$\|\mathbf{F}(\mathbf{x}_{k+1})\| = \|\mathbf{F}(\mathbf{x}_{k+1}) - \mathbf{F}(\mathbf{x}^*)\| \geq \beta \|\mathbf{x}_{k+1} - \mathbf{x}^*\|$$

(习题 8.11)。据假定，对所有 $k \geq 0$ ， $\mathbf{x}_k \neq \mathbf{x}^*$ 且 \mathbf{x}_k 恒在 D 中，故当 k 充分大后有

$$\begin{aligned} \frac{\|\mathbf{F}(\mathbf{x}_{k+1})\|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|} &\geq \frac{\beta \|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_{k+1} - \mathbf{x}^*\| + \|\mathbf{x}_k - \mathbf{x}^*\|} \\ &= \beta \frac{q_k}{1 + q_k}. \end{aligned}$$

其中 $q_k = \|\mathbf{x}_{k+1} - \mathbf{x}^*\| / \|\mathbf{x}_k - \mathbf{x}^*\|$. 据(9.3.6)知 $q_k / (1 + q_k)$ 必收敛于零, 因此 q_k 亦必收敛于零. 据定义 8.3.2 知 $\{\mathbf{x}_k\}$ 超线性收敛于 \mathbf{x}^* .

反之, 若 $\{\mathbf{x}_k\}$ 超线性收敛于 \mathbf{x}^* 且 $F(\mathbf{x}^*) = 0$, 则因

$$\frac{\|F(\mathbf{x}_{k+1})\|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|} = \frac{\|F(\mathbf{x}_{k+1}) - F(\mathbf{x}^*)\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \cdot \frac{\|\mathbf{x} - \mathbf{x}^*\|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|},$$

于是利用引理 8.4.3 及 (8.3.3)、(8.3.4) 式即知

$$\lim_{k \rightarrow \infty} \frac{\|F(\mathbf{x}_{k+1})\|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|} = 0,$$

再利用 (9.3.5) 式即可推知 (9.3.4) 式.

定理 9.3.1 设 $F: R^n \rightarrow R^n$ 满足假设 1 和假设 2, 则 Broyden 算法 9.2.1 局部、超线性收敛. 换言之, 存在正数 ε 和 δ , 当 $\|B_0 - F'(\mathbf{x}^*)\| \leq \delta$ 、 $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq \varepsilon$ 时, $\{\mathbf{x}_k\}$ 超线性收敛于 \mathbf{x}^* .

证 设 $\alpha = \|F'(\mathbf{x}^*)^{-1}\|$, 取 $\delta \leq \frac{1}{6\alpha}$, $\varepsilon \leq \frac{\delta}{2\kappa}$ 且保证 S

$(\mathbf{x}^*, \varepsilon) \subset D$. 现取 \mathbf{x}_0, B_0 满足定理要求, 于是据 Banach 引理 B_0 必为非异且

$$\|B_0^{-1}\| \leq \alpha(1 - \alpha\delta)^{-1} = \gamma.$$

又 ε 的选取方法保证了 $F(\mathbf{x}_0)$ 存在, 因此 \mathbf{x}_1 是完全确定的. 据引理 8.4.4 及定理假设应有

$$\begin{aligned} \|\mathbf{x}_1 - \mathbf{x}^*\| &= \|\mathbf{x}_0 - \mathbf{x}^* - B_0^{-1}(F(\mathbf{x}_0) - F(\mathbf{x}^*))\| \\ &\leq \|B_0^{-1}\| \|F(\mathbf{x}_0) - F(\mathbf{x}^*) - B_0(\mathbf{x}_0 - \mathbf{x}^*)\| \\ &\leq \|B_0^{-1}\| [\|F(\mathbf{x}_0) - F(\mathbf{x}^*) - F'(\mathbf{x}^*)(\mathbf{x}_0 - \mathbf{x}^*)\| \\ &\quad + \|F'(\mathbf{x}^*) - B_0\| \|\mathbf{x}_0 - \mathbf{x}^*\|] \\ &\leq \gamma \left(\frac{1}{2} \kappa \|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \delta \|\mathbf{x}_0 - \mathbf{x}^*\| \right) \end{aligned}$$

$$\begin{aligned} &\leq \gamma \left(\frac{\delta}{4} + \delta \right) \|x_0 - x^*\| \\ &< \frac{1}{2} \|x_0 - x^*\|. \end{aligned}$$

由此可见 $x_1 \in \overline{S(x^*, \varepsilon)}$, 故 $F(x_1)$ 和 B_1 存在. 再利用引理 9.3.1 即可推得

$$\begin{aligned} \|B_1 - F'(x^*)\| &\leq \|B_0 - F'(x^*)\| + \kappa \sigma(x_0, x_1) \\ &\leq \delta + \kappa \|x_0 - x^*\| \\ &\leq \delta + \frac{1}{2} \delta = \frac{3}{2} \delta < 2\delta. \end{aligned}$$

于是据 Banach 引理, B_1 为非异且

$$\|B_1^{-1}\| \leq a(1 - 2a\delta)^{-1},$$

因此 x_2 是可以确定的. 现作归纳法假定: 设 x_1, \dots, x_n 和 $B_1^{-1}, \dots, B_{n-1}^{-1}$ 均存在且当 $k \leq n$ 时有

$$\|x_k - x^*\| \leq \frac{1}{2} \|x_{k-1} - x^*\|, \quad \|B_k - F'(x^*)\| \leq \left(2 - \left(\frac{1}{2}\right)^k\right) \delta.$$

于是, 据 Banach 引理知, B_n 为非异且

$$\|B_n^{-1}\| \leq a(1 - 2a\delta)^{-1},$$

这就保证了 x_{n+1} 存在. 另一方面

$$\begin{aligned} \|x_{n+1} - x^*\| &\leq \|B_n^{-1}\| [\|F(x_n) - F(x^*) - F'(x^*)(x_n - x^*)\| \\ &\quad + \|F'(x^*) - B_n\| \|x_n - x^*\|] \\ &\leq a(1 - 2a\delta)^{-1} \left[\frac{1}{2} \kappa \|x_n - x^*\|^2 + \left(2 - \left(\frac{1}{2}\right)^n\right) \delta \|x_n - x^*\| \right] \\ &\leq a(1 - 2a\delta)^{-1} \left[\left(\frac{1}{2}\right)^n \left(\frac{\delta}{4}\right) + \left(2 - \frac{1}{2}\right)^n \delta \right] \|x_n - x^*\| \\ &< a(1 - 2a\delta)^{-1} 2\delta \|x_n - x^*\| \\ &< \frac{1}{2} \|x_n - x^*\|. \end{aligned}$$

故 $x_{n+1} \in S(x^*, \varepsilon)$. 其次

$$\begin{aligned}
\|B_{n+1} - F'(x^*)\| &\leq \|B_n - F'(x^*)\| + \kappa \sigma(x_n, x_{n+1}) \\
&\leq \left(2 - \left(\frac{1}{2}\right)^n\right) \delta + \kappa \|x_n - x^*\| \\
&\leq \left[\left(2 - \left(\frac{1}{2}\right)^n\right) + \left(\frac{1}{2}\right)^{n+1}\right] \delta \\
&= \left[2 - \left(\frac{1}{2}\right)^{n+1}\right] \delta.
\end{aligned}$$

至此归纳证明完成。

综上所述, 序列 $\{x_k\}$ 、 $\{B_k^{-1}\}$ 均存在且 $\|x_k - x^*\| \leq \left(\frac{1}{2}\right)^k \|x_0 - x^*\|$, 因此 Broyden 方法至少是线性收敛的。

以下进一步证明 Broyden 方法是超线性收敛的, 据引理 9.3.2 知只要证明 (9.3.4) 式成立即可。

通过直接计算即可证明

$$\left\|E \left(I - \frac{ss^T}{\langle s, s \rangle}\right)\right\|_F^2 = \|E\|_F^2 - \left(\frac{\|Es\|}{\|s\|}\right)^2, \quad (9.3.7)$$

其中 $E \in R^{n \times n}$ 。又由不等式 $(a^2 - \beta^2)^{1/2} \leq a - (2a)^{-1}\beta^2$ ($0 \leq \beta \leq \sqrt{2a}$) 即可推得

$$\left\|E \left(I - \frac{ss^T}{\langle s, s \rangle}\right)\right\|_F \leq \|E\|_F - (2\|E\|_F)^{-1} \left(\frac{\|Es\|}{\|s\|}\right)^2. \quad (9.3.8)$$

现令 $E = B_k - F'(x^*)$ 而 $\|E\|_F = \eta_k$, 对 (9.3.3) 利用不等式 (9.3.8) 和引理 8.4.4 即可推得

$$\eta_{k+1} \leq [1 - (2\eta_k^2)^{-1}\psi_k^2] \eta_k + \kappa \sigma_k,$$

其中

$$\psi_k = \frac{\| [B_k - F'(x^*)] s_k \|}{\|s_k\|}.$$

因为 $\{x_k\}$ 线性收敛, 故 $\sum_{k=0}^{\infty} \|x_k - x^*\| < +\infty$. 另一方面, 由引理

9.3.1知, $\eta_{k+1} \leq \eta_k + \kappa \sigma_k$, 据此不等式即可推知数列 $\{\eta_k\}$ 有界.

设 η 为它的某一上界, 则对所有 k 均成立

$$(2\eta^2)^{-1} \psi_k^2 \leq \eta_k - \eta_{k+1} + \kappa \sigma_k,$$

于是

$$\begin{aligned} (2\eta^2)^{-1} \sum_{k=0}^{\infty} \psi_k^2 &\leq \eta_0 + \kappa \sum_{k=0}^{\infty} \sigma_k \\ &\leq \delta + \kappa \sum_{k=0}^{\infty} \left(\frac{1}{2}\right)^k \|x_0 - x^*\| \leq \delta + 2\kappa\epsilon. \end{aligned}$$

由此可见 $\{\psi_k\}$ 收敛于零, 亦即(9.3.4)式成立.

第九章 习 题

9.1 设 $a, b \in R^n$, 试证: $\|ab^T\|_2 = \|a\|_2 \|b\|_2 = \|ab^T\|_F$.

9.2 设 $a, b \in R^n$ 且 $a^T b = 1$, 试证:

i) $\|I - ab^T\|_2 = \|a\|_2 \|b\|_2$;

ii) $\|I - ab^T\|_2 = 1$ 的充要条件为 a 和 b 共线.

9.3 设 $A, B \in R^{n \times n}$, 试证:

$$\|AB\|_F \leq \min\{\|A\|_2 \|B\|_F, \|A\|_F \|B\|_2\}.$$

9.4 设给定非零向量 $z, y \in R^n, A, B \in R^{n \times n}$, 试验证:

$$f(A) = \|A - B\|_F$$

为凸集 $N = \{A \mid Az = y\}$ 上严格的凸函数.

9.5 设 $A \in R^{n \times n}$ 为非异, $u, v \in R^n$, 试证:

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A u},$$

其中 $1 + v^T A u \neq 0$ (Sherman-Morrison公式).

9.6 设 $B \in R^{n \times n}, S \in R^{k \times k}$ 而 $U, V \in R^{n \times k}$. 假定 B, S 和 $V^T B^{-1} U - S^{-1}$ 为非异, 试证:

$$(B + USV^T)^{-1} = B^{-1} - B^{-1} U T V^T B^{-1}$$

其中 $S^{-1} + T^{-1} = V^T B^{-1} U$ (Sherman-Morrison-Woodbury公式).

9.7 设 $F: R^n \rightarrow R^n$ 满足引理9.3.2的假设, 如果对某 $x_0 \in D$ 由

$$x_{k+1} = x_k - \lambda_k B_k^{-1} F(x_k) \quad (k=0, 1, \dots)$$

生成的迭代序列 $\{x_k\}$ 是完全确定的且 $\{x_k\} \subset D$ 收敛于 x^* , 又假定(9.3.4)式成立, 那末 $\{x_k\}$ 超线性收敛于 x^* , 当且仅当 $\{\lambda_k\}$ 收敛于1.

9.8 利用算法9.2.1计算例9.2.1.

9.9 试证等式(9.3.7).

参 考 文 献

1. Dennis, J.E. & Moré, J.J. (1977), "Quasi Newton Methods, Motivation and Theory", SIAM Review, Vol.19, No.1, January.
2. Broyden, G.G. Dennis, J.E. & Moré, J.J. (1973), "On the local and superlinear Convergence of quasi-Newton methods", Ibid., 12, pp.223—246.
3. Dennis, J.E. & More, J. J. (1974), "A Characterization of superlinear Convergence and its application to quasi-Newton methods", Math Comp, 28, pp.549—560.
4. 王德人 非线性方程组解法与最优化方法.
5. Ralston, A. & Rabinowitz, P. (1978), A First Course in Numerical Analysis.
6. Gill, P.E. & Murray, W. (1972), "Quasi-Newton methods for unconstrained minimization", J.Inst. Math. Appl., 9, pp.91—108.
7. Stewart, G.w., (1973). Introduction to Matrix Computation. New York, Academic Press.